

### ***Remarks***

Reconsideration of this Application is respectfully requested.

Claims 38-42, 45-49 and 51-57 are pending in the application, with claims 38 and 51 being the independent claims. Claims 1-37, 43, 44 and 50 are cancelled. Claims 38, 47, 48, 49 and 51 have been amended. Support for the amendments to claims 38 and 51 claims may be found throughout the specification. Support for the amendment to claims 47-49 may be found throughout the specification, e.g., in paragraph [0072] on page 23.

Based on the following remarks, Applicants respectfully request that the Examiner reconsider all outstanding rejections and that they be withdrawn.

### ***Examiner Interview***

On June 6, 2007, the Examiner initiated a telephonic interview with Applicants representatives. The Examiner suggested additional amendments to the claims to place the claims in condition for allowance. Applicants thank the Examiner for suggesting the claim amendments and have amended the claims in accordance with the Examiner's suggestions.

The Examiner further requested that Applicants submit copies of Declarations Under 37 C.F.R. § 1.132 by En Li, Ph.D., an inventor, and Kenneth D. Bloch, M.D. that were filed in parent Appl. No. 09/720,086 in support of Applicants' priority claim. Copies of these documents are submitted herewith.

Accordingly, Applicants respectfully request that the Examiner reconsider and withdraw all outstanding objections and rejections.


***Conclusion***

All of the stated grounds of objection and rejection have been properly traversed, accommodated, or rendered moot. Applicants therefore respectfully request that the Examiner reconsider all presently outstanding objections and rejections and that they be withdrawn. Applicants believe that a full and complete reply has been made to the outstanding Office Action and, as such, the present application is in condition for allowance. If the Examiner believes, for any reason, that personal communication will expedite prosecution of this application, the Examiner is invited to telephone the undersigned directly at (202) 772-8658.

Prompt and favorable consideration of this Reply is respectfully requested.

Respectfully submitted,

STERNE, KESSLER, GOLDSTEIN & FOX P.L.L.C.

  
Daniel J. Nevriy  
Agent for Applicants  
Registration No. 59,118

Date: June 7, 2007

1100 New York Avenue, N.W.  
Washington, D.C. 20005-3934  
(202) 371-2600



IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re application of:

Li *et al.*

Appl. No. 09/720,086

102(e): July 23, 2001

For: **De Novo DNA Cytosine  
Methyltransferase Genes,  
Polypeptides and Uses Thereof**

Confirmation No.: 6968

Art Unit: 1642

Examiner: Harris, A. M.

Atty. Docket: 0609.4560002/KRM/DJN

**Declaration Under 37 C.F.R. § 1.132 of En Li, Ph.D.**

Commissioner for Patents  
P.O. Box 1450  
Alexandria, VA 22313-1450

Sir:

I, the undersigned, En Li, Ph.D., residing at 45 Hinckley Road, Newton, Massachusetts, 02168, declare and state as follows:

1. I am a co-inventor of the above-captioned patent application.
2. I am currently employed by Novartis Institutes for Biomedical Research as Vice President & Global Head, Animal Models of Disease and Epigenetics Program. Prior to my current employment, I was an Associate Professor of Medicine at Harvard Medical School and directed a laboratory in the Cardiovascular Research Center at the Massachusetts General Hospital from January 1993 to April 2003, where I conducted and supervised research in the field of mouse genetics and developmental biology.
3. A current *curriculum vitae* is appended hereto as EXHIBIT A.
4. I have reviewed the above-captioned patent application and the Office Action dated June 6, 2005. I have also reviewed the sequence listing as filed and the sequence listing as amended on July 23, 2001. I have also reviewed the claims of the captioned patent application.
5. I have been informed that the Examiner has not granted priority to the earlier filed patent applications because there is insufficient proof that the coding regions of currently amended SEQ ID NOS:1 and 2 are the same as those listed in the priority documents, viz., the mouse Dnmt3a and Dnmt3b cDNA clones encoding the coding regions of SEQ ID NOS:1 and 2, respectively, that were deposited with the American Type Culture Collection (ATCC), 10801 University Boulevard, Manassas, Virginia 20110-2209, USA. Sequences harboring the coding regions of SEQ ID NOS: 1 and 2, respectively, were

- 2 -

Li et al.  
Appl. No. 09/720,086

deposited with the ATCC on June 16, 1998, and assigned ATCC Deposit Nos. 209933 and 209934, respectively. The deposit date of June 16, 1998 was prior to the filing date of the first provisional application, App. No. 60/090,906, filed June 25, 1998, the benefit of which is claimed. The '906 application includes the sequence information and references the deposits of the sequenced material on page 15, lines 26, through page 16, line 2, of the specification.

6. In November 2004, Applicants had samples withdrawn of the mouse Dnmt3a and Dnmt3b cDNA clones contained within ATCC Deposit Nos. 209933 and 209934, respectively. At the Applicants request, Kenneth D. Bloch, M.D., a faculty member in the Cardiovascular Research Center at the Massachusetts General Hospital and an experienced DNA sequencer, sequenced nucleotides that spanned the coding regions of mouse Dnmt3a and Dnmt3b in the deposited cDNA. A nucleotide alignment that spans the coding regions of sequenced mouse Dnmt3a cDNA clone contained in ATCC Deposit No. 209933 and currently amended **SEQ ID NO:1** is shown in EXHIBIT B. A nucleotide alignment that spans the coding regions of sequenced mouse Dnmt3b cDNA clone contained in ATCC Deposit No. 209934 and currently amended **SEQ ID NO:2** is shown in EXHIBIT C.

7. The amendment to the sequence listing, which was filed on July 23, 2001, corrected six nucleotides in the coding sequence of **SEQ ID NO:1** (see the bolded nucleotides at positions 516, 843, 1036, 1110, 1116 and 1726 in EXHIBIT B) and two nucleotides in the coding sequence of **SEQ ID NO:2** (see the bolded nucleotides at positions 918 and 920 in EXHIBIT C).

8. The deposited clones recited in ¶¶5 and 6, above (i.e., ATCC Deposit Nos. 209933 and 209934) are the same as the deposited clones recited in the above-captioned application. The coding sequence of ATCC Deposit No. 209933 is currently believed to be the same as the coding sequence of currently amended **SEQ ID NO:1**. The coding sequence of ATCC Deposit No. 209934 is currently believed to be the same as the coding sequence of currently amended **SEQ ID NO:2**.

9. It is well known that sequencing errors are a common problem in Molecular Biology. See, e.g., Peter Richterich, Estimation of Errors in 'Raw' DNA Sequences: A Validation Study, 8 *Genome Research* 251-59 (1998)(EXHIBIT D). I believe that one skilled in the art would have sequenced the deposited material and recognized the sequencing errors in the coding region. I believe that the correct mouse Dnmt3a and Dnmt3b coding sequences are inherent to the ATCC deposited clones, ATCC Deposit Nos. 209933 and 209934, respectively, which were deposited prior to the filing of App. No. 60/090,906, filed June 25, 1998, the benefit of which is claimed.

10. Accordingly, based on the above, I believe that Applicants are entitled to the June 25, 1998 filing date for the coding sequences of mouse Dnmt3a and Dnmt3b contained within ATCC Deposit Nos. 209933 and 209934, respectively.

- 3 -

Li et al.  
Appl. No. 09/720,086

11. I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the present patent application or any patent issued thereon.

Respectfully submitted,



En Li, Ph.D.

Date: 11/7/05

304890v1

# EXHIBIT A

# CURRICULUM VITAE

**En Li, Ph.D.**

Vice President & Global Head  
Animal Models of Disease  
Novartis Institute for Biomedical Research, Inc.  
250 Mass Ave. Cambridge, MA 02139  
Tel: 617-871-7072  
Fax: 617-871-7263  
Email: en.li@novartis.com

## **Education:**

1984 B.Sc. in Biochemistry. Peking University, Beijing, China  
1992 Ph.D. Massachusetts Institute of Technology (Biology)  
(Advisor: Prof. Rudolf Jaenisch)

## **Professional Experience:**

1993-1996 Principal Investigator, Cardiovascular Research Center, Massachusetts General Hospital  
Instructor, Department of Medicine, Harvard Medical School  
1996-2000 Assistant Professor, Department of Medicine, Harvard Medical School  
2000-2003 Associate Professor, Department of Medicine, Harvard Medical School  
2001-2003 Guest Professor, Beijing University, Health Science Center  
2003- VP & Global Head, Models of Disease Center, Epigenetics Program  
Novartis Institute for Biomedical Research.

## ***Review and editorial board***

1999- External Grant Reviewer, Human Frontier Science Program  
1999- External Grant Reviewer, NIH, NICHD  
1999- Member of the Advisory Board. Journal of Biochemistry  
2000- External Grant Reviewer, NIH, NIA  
2000- External Grant Reviewer, NSF  
2000- Mail Reviewer, Wellcome Trust  
2003- Member of the Advisory Board. China Science Reports  
2004 External Grant Reviewer, Chinese Natural Science Foundation.  
  
1993- Ad Hoc Reviewer for the following journals  
Nature, Science, Cell, Nat Genet, Genes Dev, Trends Genet, Development, Mol Cell Biol, PNAS, Human Mol. Genet., Dev. Biol., Mech. Dev., J. Cell Biol., Dev. Dyn., Gene, Genomics, Nucl Acid Res, Mammalian Genome, etc.

## ***Professional Societies***

1990- American Association for the Advancement of Science, Member

1994- DNA Methylation Society, Member  
1999- Ray Wu Society, Member

**Invited Presentations (Since 2003-)**

2003 Invited speaker, Gordon Research Conference - Cancer Genetics and Epigenetics  
2003 Session chair, Keystone Symposium – Chromatin (Big Sky, Montana)  
2003 Invited speaker, Annual HUGO meeting at Cancun, Mexico  
2003 Invited speaker, Gordon Research Conference – Epigenetics  
2004 Invited speaker, the 2<sup>nd</sup> Annual CDB symposium, Kobe, Japan  
2004 Invited speaker, 2<sup>nd</sup> Weissenburg Symposium on DNA methylation – an important genetics signals. Weissenburg, Germany  
2004 Session Chair, on genomic imprinting. 10<sup>th</sup> SCBA International Symposiums, Beijing, China  
2004 Invited speaker, Genomic imprinting workshop in Montpellier, France  
2005 Vice Chair, Gordon Research Conference 'Cancer Genetics and Epigenetics'



## **Publication:**

1. Zijlstra M, Li E, Sajjadi F, Subramani S, Jaenisch R. Germ-line transmission of a disrupted b2-microglobulin gene produced by homologous recombination in embryonic stem cells. *Nature* 1989; 342: 435-8.
2. Lee K, Li E, Huber J, Landis S, Sharpe A, Chao MV, Jaenisch R. Targeted mutation of the p75 low affinity NGF receptor gene leads to deficits in the peripheral sensory nervous system. *Cell* 1992; 69: 737-49.
3. Li E, Bestor TH, Jaenisch R. Targeted mutation of the DNA methyltransferase gene results in embryonic lethality. *Cell* 1992; 69: 915-26.
4. Li E, Sucov HM, Lee K, Evans RM, Jaenisch R. Normal development and growth of mice carrying a targeted disruption of the  $\alpha 1$  retinoic acid receptor gene. *Proc Natl Acad Sci USA* 1993; 90: 1590-4.
5. Li E, Beard C, Jaenisch R. Role for DNA methylation in genomic imprinting. *Nature* 1993; 366: 362-5.
6. Jüttermann R, Li E, Jaenisch R. The toxicity of 5-aza-2'-deoxycytidine to mammalian cells is mediated primarily by DNA methyltransferase rather than DNA demethylation. *Proc Natl Acad Sci USA* 1994; 91: 11797-11801.
7. Laird PW, Jackson-Grusby L, Fazeli A, Dickinson W, Jung E, Li E, Weinberg RA, Jaenisch R. Suppression of intestinal neoplasia by DNA hypomethylation. *Cell* 1995; 81:197-205.
8. Shinoda K, Lei H, Yoshii H, Nomura M, Nagano M, Shiba H, Sasaki H, Osawa Y, Ninomiya Y, Niwa O, Morohashi K, Li E. Developmental defects of the ventromedial hypothalamic nucleus and pituitary gonadotroph in the Ftz-F1 disrupted mice. *Dev Dyn* 1995; 204: 22-9.
9. Beard C, Li E, Jaenisch J. Loss of methylation activates Xist in somatic but not in embryonic cells. *Genes Dev* 1995; 9: 2325-2334.
10. Nan X, Tate P, Li E, Bird A. DNA methylation specifies chromosomal localization of MeCP2. *Mol Cell Biol* 1996; 16: 414-421.
11. Zhang W, Zimmer G, Chen J, Ladd D, Li E, Alt FW, Wiederrecht G, Cryan J, O'Neill EA, Seidman CE, Abbas AK, Seidman JG. T cell responses in calcineurin  $\text{A}\alpha$ -deficient mice. *J Exp Med* 1996; 183: 413-420.
12. Tucker KL, Beard C, Dausman J, Jackson-Grusby L, Laird PW, Lei H, Li E, Jaenisch R. Germ-line passage is required for establishment of methylation and expression patterns of imprinted but not of nonimprinted genes. *Genes Dev* 1996; 10: 1008-1020.
13. Trasler JM, Trasler DG, Bestor TH, Li E, Ghibu F. Nuclear localization and expression of DNA methyltransferase in embryos is temporally correlated with postimplantation increase in DNA methylation. *Dev Dyn* 1996; 206: 239-247.

14. Harbers K, Müller U, Grams A, Li E, Jaenisch R, Franz T. Provirus integration into a gene encoding a ubiquitin-conjugating enzyme results in a placental defect and embryonic lethality. *Proc Natl Acad Sci USA* 1996; 93: 12412-7.
15. Nakamuta M, Chang B, Zsigmond E, Kobayashi K, Lei H, Ishida BY, Oka K, Li E, Chan L. Complete phenotypic characterization of apobec-1 knockout mice with a wild-type genetic background and a human apolipoprotein B transgenic background, and restoration of apolipoprotein B mRNA editing by somatic gene transfer of apobec-1. *J Biol Chem* 1996; 271: 25981-8.
16. Lei H, Oh SP, Okano M, Juttermann R, Goss KA, Jaenisch R, Li E. De novo DNA cytosine methyltransferase activities in mouse embryonic stem cells. *Development* 1996; 122: 3195 - 3205.
17. Chae T, Kwon YT, Bronson R, Dikkes P, Li E, Tsai L-H. Mice lacking p35, a neuronal specific activator of Cdk5, display cortical lamination defects, seizures, and adult lethality. *Neuron* 1997; 18: 29-42.
18. Simon AM, Goodneough DA, Li E, Paul DL. Female infertility in mice lacking connexin 37. *Nature* 1997; 385: 525-9.
19. Oh SP, Li E. The signaling pathway mediated by type IIB activin receptor controls axial patterning and lateral asymmetry in the mouse. *Genes Dev* 1997; 11: 1812-1826.
20. Gong X, Li E, Klier G, Huang Q, Wu Y, Lei H, Kumar NM, Horwitz J, Gilula NB. Disruption of  $\alpha 3$  connexin gene leads to protein degradation and cataractogenesis in mice. *Cell* 1997; 91: 833-843.
21. Bonventre JV, Huang Z, Taheri MR, O'Leary E, Li E, Moskowitz MA, Sapirstein, A. Cytosolic phospholipase A2 deficient mice have abnormal fertility and are protected against focal ischemic injury to the brain. *Nature* 1997; 390: 622-5.
22. Pradhan S, Talbot D, Sha M, Benner J, Hornstra L, Li E, Jaenisch R, Roberts RJ. Baculovirus-mediated expression and characterization of the full-length murine DNA methyltransferase. *Nucl Acids Res* 1997; 25: 4666-4673.
23. Wang S, Miura M, Jung Y-K, Zhu H, Li E, Yuan J. Murine caspase-11, an ICE-interacting protease is essential for the activation of ICE. *Cell* 1998; 92: 501-9.
24. Gu Z, Nomura M, Simpson BB, Lei H, van den Eijnden-van Raaij AJM, Donahoe PK, Li E. The type I activin receptor ActR-IB is required for egg cylinder organization and gastrulation in the Mouse. *Genes Dev* 1998; 12: 844-857.
25. Okano M, Xie S, Li E. Dnmt2 is not required for de novo and maintenance methylation of viral DNA in embryonic stem cells. *Nucl Acids Res* 1998; 26: 2536-2540.
26. Bergeron L, Perez GI, MacDonald G, Shi L, Sun Y, Jurisicova A, Varmuza S, Latham K.E Flaws JA, Salter JC, Hara H, Moskowitz MA, Li E, Greenberg A, Tilly JL, Yuan J. Defects in regulation of apoptosis in caspase-2-deficient mice. *Genes Dev* 1998; 12: 1304-1314.

27. Yao TP, Oh SP, Fuchs M, Zhou ND, Ch'ng LE, Newsome D, Bronson RT, Li E, Livingston DM, Eckner R. Gene dosage-dependent embryonic development and proliferation defects in mice lacking the transcriptional integrator p300. *Cell* 1998; 93: 361-372.
28. Chung UI, Lanske B, Lee K, Li E, Kronenberg, H. The parathyroid hormone/parathyroid hormone-related peptide receptor coordinates endochondral bone development by directly controlling chondrocyte differentiation. *Proc Natl Acad Sci USA* 1998; 95: 13030-5.
29. Nomura M, Li E. Roles for Smad2 in mesoderm formation, left-right patterning, and craniofacial development in mice. *Nature* 1998; 393: 786-790.
30. Okano M, Xie S, Li E. Cloning and characterization of a family of novel mammalian DNA (cytosine-5) methyltransferases. *Nat Genet* 1998; 19: 219-220.
31. Morohashi K, Tsuboi-Asai H, Matsushita S, Suda M, Nakashima M, Sasano H, Hataba Y, Li CL, Fukata J, Irie J, Watanabe T, Nagura H, Li E. Structural and functional abnormalities in the spleen of an mFtz-F1 gene-disrupted mouse. *Blood* 1999; 93: 1586-1594.
32. Gu Z, Reynolds EM, Song J, Lei H, Feijen A, Yu L, He W, MacLaughlin DT, van den Eijnden-van Raaij J, Donahoe PK, Li, E. The type I serine/threonine kinase receptor ActRIA (ALK2) is required for gastrulation of the mouse embryo. *Development* 1999; 126: 2551-2561.
33. Verheijen MH, Karperien M, Chung U, van Wijnen M, Heystek H, Hendriks JA, Veltmaat JM, Lanske B, Li E, Lowik CW, de Laat SW, Kronenberg HM, Defize LH. Parathyroid hormone-related peptide (PTHrP) induces parietal endoderm formation exclusively via the type I PTH/PTHrP receptor. *Mech Dev* 1999; 81: 151-161.
34. Song J, Oh SP, Schrewe H, Nomura M, Lei H, Okano M, Gridley T, Li E. The type II activin receptors are essential for egg cylinder growth, gastrulation and rostral head development in mice. *Dev Biol* 1999; 213: 157-169.
35. Xie S, Wang Z, Okano M, Nogami M, Li Y, He W, Okumura K, Li E. Cloning, expression, and chromosome locations of the human DNMT3 gene family. *Gene* 1999; 236: 87-95.
36. Donahoe MA, Zhang X-L, McGinnis L, Biggers J, Li E, Shi Y. Targeted disruption of mouse Yin Yang 1 transcription factor results in peri-implantation lethality. *Mol Cell Biol* 1999; 19: 7237-7244.
37. Okano M, Bell DW, Haber DA, Li E. DNA methyltransferases Dnmt3a and Dnmt3b are essential for de novo methylation and mammalian development. *Cell* 1999; 99: 247-257.
38. Okano M, Takebayashi S, Okumura K, Li E. Assignment of cytosine-5 DNA methyltransferases Dnmt3a and Dnmt3b to mouse chromosome bands 12A2-A3 and 2H1 by in situ hybridization. *Cytogenet Cell Genet* 1999; 86: 333-334.
39. Li YP, Chen W, Liang Y, Li E, Stashenko P. OC-116Kda-deficient mice exhibit severe osteopetrosis due to loss of osteoclast-mediated extracellular acidification. *Nat Genet* 1999; 23: 452-456.

40. Nakagawa T, Zhu H, Morishima N, Li E, Xu J, Yankner BA, Yuan J. Caspase-12 mediates endoplasmic-reticulum-specific apoptosis and cytotoxicity by amyloid-beta. *Nature* 2000; 403:98-103.
41. Oh, SP., Seki, T., Goss, KA., Imamura, T., Yi, Y., Donahoe, PK., Li, L., Miyazono, K., ten Dijke, P., Kim, S., and Li, E. Activin Receptor-Like Kinase 1 (ALK1) modulates TGF- $\beta$ 1 signaling in the regulation of angiogenesis. *Proc Natl Acad Sci USA* 2000 (in press).
42. Sado, T., Fenner, M. H., Tan, S.-S., Tam, P., Shioda, T., and Li, E. (2000). X inactivation in the mouse embryo deficient for *Dnmt1*: distinct effect of hypomethylation on imprinted and random X inactivation. *Dev. Biol.* 15: 294-303.
43. Kim, S.K., Hebrok, M., Li, E., Oh, S.P., Schrewe, H., Harmon, E.B., Lee, J.S., and Melton, D.A. (2000). Activin receptor patterning of foregut organogenesis. *Genes Dev* 14:1866-1871.
44. Pannell D, Osborne CS, Yao S, Sukonnik T, Pasceri P, Karauskakis A, Okano M, Li E, Lipshitz HD, and Ellis J (2000) Retrovirus vector silencing is de novo methylase independent and marked by a repressive histone code. *EMBO J.* 19:5884-5894.
45. Sado, T., Wang, Z., Sasaki, H., and Li, E. (2001). Regulation of Imprinted X-Inactivation in Mice by *Tsix*. *Development* 128: 1275-1286.
46. Ferguson, C.A., Tucker, A.S., Heikinheimo, K., Nomura, M., Oh, P., Li, E., Sharpe, PT. (2001) The role of effectors of the activin signalling pathway, activin receptors IIA and IIB, and Smad2, in patterning of tooth. *Development* 128:4605-4613.
47. Ko, J., Humbert, S., Bronson, R.T., Takahashi, S., Kulkarni, A.B., Li, E., Tsai, L.H. (2001) p35 and p39 are essential for cyclin-dependent kinase 5 function during neurodevelopment. *J Neurosci* 21:6758-6771.
48. Jiang P, Song J, Gu G, Slonimsky E, Li E, Rosenthal N. (2002). Targeted deletion of the *MLC1f/3f* downstream enhancer results in precocious MLC expression and mesoderm ablation. *Dev Biol* 243:281-293.
49. Liang G, Chan MF, Tomigahara Y, Tsai YC, Gonzales FA, Li E, Laird PW, Jones PA. (2002). Cooperativity between DNA Methyltransferases in the Maintenance Methylation of Repetitive Elements. *Mol Cell Biol* 22(2):480-491.
50. Hata, H., Okano, M., Lei, H., and Li, E. (2002) *Dnmt3L* cooperates with the *Dnmt3* family de novo DNA methyltransferases to establish maternal imprinting in mice. *Development* 129: 1983-1993.
51. Sakuma R, Ohnishi Yi Y, Meno C, Fujii H, Juan H, Takeuchi J, Ogura T, Li E, Miyazono K, Hamada H. (2002) Inhibition of Nodal signalling by Lefty mediated through interaction with common receptors and efficient diffusion. *Genes Cells* 7:401-412.
52. Miura K, Kishino T, Li E, Webber H, Dikkes P, Holmes GL, Wagstaff J. (2002) Neurobehavioral and electroencephalographic abnormalities in *ube3a* maternal-deficient mice. *Neurobiol Dis* 9:149-159.

53. Dodge, J. Ramsahoye, B.H., Wo, Z.G, Okano, M., and Li, E. (2002) *De novo* methylation of MMLV provirus in embryonic stem cells: CpG versus non-CpG methylation. *Gene* 289:41-48.
54. Oh, S.P., and Li, E. (2002). Gene-dosage sensitive genetic interactions between *inversus viscerum* (*iv*), *nodal*, and activin type IIB receptor (*ActRIIB*) genes in asymmetrical patterning of the visceral organs along the left-right axis. *Dev. Dyn.* 224:279-290.
55. Chen T, Ueda Y, Xie S, Li E. (2002). A novel Dnmt3a isoform produced from an alternative promoter localizes to euchromatin and its expression correlates with active de novo methylation. *J Biol Chem.* 277(41):38746-54.
56. Oh, SP, Yeo, CY, Lee, Y, Schrewe, H, Whitman, M, and Li, E (2002). Activin type IIA and type IIB receptors mediate Gdf11 signaling in axial vertebral patterning. *Genes Dev.* 16: 2749-2754.
57. Sado T, Li E, Sasaki H. Effect of TSIX disruption on XIST expression in male ES cells. *Cytogenet Genome Res.* 2002;99:115-118.
58. Trasler J, Deng L, Melnyk S, Pogribny I, Hiou-Tim F, Sibani S, Oakes C, Li E, James SJ, Rozen R. Impact of Dnmt1 deficiency, with and without low folate diets, on tumor numbers and DNA methylation in Min mice. *Carcinogenesis.* 2003 Jan;24(1):39-45.
59. Welte T, Zhang SS, Wang T, Zhang Z, Hesslein DG, Yin Z, Kano A, Iwamoto Y, Li E, Craft JE, Bothwell AL, Fikrig E, Koni PA, Flavell RA, Fu XY(2003). STAT3 deletion during hematopoiesis causes Crohn's disease-like pathogenesis and lethality: a critical role of STAT3 in innate immunity. *Proc Natl Acad Sci U S A.* 2003 Feb 18;100(4):1879-84.
60. Wei Chen, Yuqiong Liang, Wenjie Deng, Ken Shimizu, Amir M. Ashique, En Li, and Yi-Ping Li (2003). The zinc-finger protein CNBP is required for forebrain formation in the mouse. *Development* 130: 1367-1379.
61. Joost Gribnau, Konrad Hochedlinger, Ken Hata, En Li, and Rudolf Jaenisch (2003). Asynchronous replication timing of imprinted loci is independent of DNA methylation, but consistent with differential subnuclear localization. *Genes Dev* 17:759-773.
62. Lehnertz B, Ueda Y, Derijck AA, Braunschweig U, Perez-Burgos L, Kubicek S, Chen T, Li E, Jenuwein T, Peters AH. (2003) Suv39h-mediated histone h3 lysine 9 methylation directs DNA methylation to major satellite repeats at pericentric heterochromatin. *Curr Biol.* 13:1192-1200.
63. Genomic organization and promoter analysis of the Dnmt3b gene. (2003) Ishida C, Ura K, Hirao A, Sasaki H, Toyoda A, Sakaki Y, Niwa H, Li E, Kaneda Y. *Gene.* 310:151-9.
64. Taiping Chen, Yoshihide Ueda, Jonathan E. Dodge, Zhenjuan Wang, and En Li (2003) Establishment and Maintenance of Genomic Methylation Patterns in Mouse Embryonic Stem Cells by Dnmt3a and Dnmt3b. *Mol. Cell Biol.*, 23:5594-5605.

65. Jacoby JJ, Kalinowski A, Liu MG, Zhang SS, Gao Q, Chai GX, Ji L, Iwamoto Y, Li E, Schneider M, Russell KS, Fu XY. Cardiomyocyte-restricted knockout of STAT3 results in higher sensitivity to inflammation, cardiac fibrosis, and heart failure with advanced age. *Proc Natl Acad Sci U S A*. 2003. 100:12929-34.
66. Kelly TL, Li E, Trasler JM. 5-aza-2'-deoxycytidine induces alterations in murine spermatogenesis and pregnancy outcome. *J Androl*. 2003. 24:822-830.
67. Mund C, Musch T, Strodicke M, Assmann B, Li E, Lyko F. Comparative analysis of DNA methylation patterns in transgenic *Drosophila* overexpressing mouse DNA methyltransferases. *Biochem J*. 2004. 378:763-768.
68. Sado T, Okano M, Li E, Sasaki H. De novo DNA methylation is dispensable for the initiation and propagation of X chromosome inactivation. *Development*. 2004. 131:975-982.
69. Dodge JE, Kang YK, Beppu H, Lei H, Li E. Histone H3-K9 methyltransferase ESET is essential for early development. *Mol Cell Biol*. 2004. 24:2478-86.
70. Kaneda M, Okano M, Hata K, Sado T, Tsujimoto N, Li E, Sasaki H. Essential role for de novo DNA methyltransferase Dnmt3a in paternal and maternal imprinting. *Nature*. 2004. 429(6994):900-3.
71. Beppu H, Ichinose F, Kawai N, Jones RC, Yu PB, Zapol WM, Miyazono K, Li E, Bloch KD. BMPR-II heterozygous mice have mild pulmonary hypertension and an impaired pulmonary vascular remodeling response to prolonged hypoxia. *Am J Physiol Lung Cell Mol Physiol*. 2004 Jul 30
72. de Sousa Lopes SM, Roelen BA, Monteiro RM, Emmens R, Lin HY, Li E, Lawson KA, Mummery CL. BMP signaling mediated by ALK2 in the visceral endoderm is necessary for the generation of primordial germ cells in the mouse embryo. *Genes Dev*. 2004 Aug 1;18(15):1838-49.

73. Hattori N, Abe T, Hattori N, Suzuki M, Matsuyama T, Yoshida S, Li E, Shiota K. Preference of DNA methyltransferases for CpG islands in mouse embryonic stem cells. *Genome Res.* 2004. 14(9):1733-40.
74. Chen T, Tsujimoto N, Li E. The PWWP Domain of Dnmt3a and Dnmt3b Is Required for Directing DNA Methylation to the Major Satellite Repeats at Pericentric Heterochromatin. *Mol Cell Biol.* 2004. 24:9048-58.
75. Fang J, Chen T, Chadwick B, Li E, Zhang Y. Ring1b-mediated H2A ubiquitination associates with inactive X chromosomes and is involved in initiation of X inactivation. *J Biol Chem.* 2004 Dec 17;279(51):52812-5.
76. Beppu H, Ichinose F, Kawai N, Jones RC, Yu PB, Zapol WM, Miyazono K, Li E, Bloch KD. BMPR-II heterozygous mice have mild pulmonary hypertension and an impaired pulmonary vascular remodeling response to prolonged hypoxia. *Am J Physiol Lung Cell Mol Physiol.* 2004 Dec;287(6):L1241-7.
77. Park S, Lee YJ, Lee HJ, Seki T, Hong KH, Park J, Beppu H, Lim IK, Yoon JW, Li E, Kim SJ, Oh SP. B-cell translocation gene 2 (Btg2) regulates vertebral patterning by modulating bone morphogenetic protein/smad signaling. *Mol Cell Biol.* 2004 Dec;24(23):10256-62.
78. Feng J, Chang H, Li E, Fan G. Dynamic expression of de novo DNA methyltransferases Dnmt3a and Dnmt3b in the central nervous system. *J Neurosci Res.* 2005 Mar 15;79(6):734-46.
79. Beppu H, Lei H, Bloch KD, Li E. Generation of a floxed allele of the mouse BMP type II receptor gene. *Genesis.* 2005 Feb 25;41(3):133-137
80. Yu PB, Beppu H, Kawai N, Li E, Bloch KD. Bone morphogenetic protein (BMP) type II receptor deletion reveals BMP ligand-specific gain of signaling in pulmonary artery smooth muscle cells. *J Biol Chem.* 2005. 280:24443-50.
81. Hopfer U, Fukai N, Hopfer H, Wolf G, Joyce N, Li E, Olsen BR. Targeted disruption of Col8a1 and Col8a2 genes in mice leads to anterior segment abnormalities in the eye. *FASEB J.* 2005. 19:1232-44.
82. Rodic N, Oka M, Hamazaki T, Murawski MR, Jorgensen M, Maatouk DM, Resnick JL, Li E, Terada N. DNA methylation is required for silencing of Ant4, an adenine nucleotide translocase selectively expressed in mouse embryonic stem cells and germ cells. *Stem Cells.* 2005 Jul 28; [Epub ahead of print]

#### **Proceedings of Meetings**

1. Li E, Beard C, Foster A, Bestor TH, Jaenisch R. DNA methylation, genomic imprinting, and mammalian development. *Proceedings of the Cold Spring Harbor Symposium on Quantitative Biology*; 1993; 58: 297-305.
2. Muragaki Y, Timmons S, Griffith CM, Oh SP, Li E, Fukai N, Fadel B, Quertermous T,

Olsen B. Tissue-specific expression of three alternative variants of *Col18a1* and their localization in basement membrane zones. Proceedings of the 7th International Symposium on Basement Membranes, Bethesda, 1995.

3. Okano M, Li E. (2002). Genetic analyses of DNA methyltransferase genes in mouse model system. J Nutr. 132:2462S-5S.

#### **Reviews, Book Chapters, and Editorials**

1. Jaenisch R, Beard C, Li, E. DNA methylation and mammalian development. In: Ohlsson R, Hall K, Ritzen M, editors. Genomic imprinting: Causes and consequences. Cambridge, UK: Cambridge University Press; 1995. p. 118.

2. Li, E. Role of DNA methylation in mammalian development. In: Reik W, Sorani A, editors. Genomic Imprinting - Frontiers in Molecular Biology. Volume 18. Oxford, UK: Oxford University Press; 1997. p. 1-20.

3. Li, E. The mojo of methylation [News & Views]. Nat Genet 1999; 23: 5-6.

4. Li, E and Jaenisch, R. DNA methylation and methyltransferases. In: Ehrlich M, editor. DNA alterations in cancer: Genetic and epigenetic changes. PP. 351-365. Natick: BioTechniques Books; 2000.

5. Li E. Chromatin modification and epigenetic reprogramming in mammalian development. 2002. Nature Rev. Genet. 3, 662-673.

6. Chen T, Li E. Structure and function of eukaryotic DNA methyltransferases. 2004. Curr Top Dev Biol. 60:55-89.



# EXHIBIT B

## EXHIBIT B

**Alignment spanning the coding region of mouse Dnmt3a sequence from ATCC  
Deposit No. 209933 (top) and currently amended SEQ ID NO:1 (bottom)<sup>1</sup>**

```

atgccctccagcggccccggggacaccagcagctcctctctggagcgggaggatgatcga
|||||
atgccctccagcggccccggggacaccagcagctcctctctggagcgggaggatgatcga 276

aaggaaggagaggaacaggaggagaaccgtggcaaggaagagcgccaggagcccagcgcc
|||||
aaggaaggagaggaacaggaggagaaccgtggcaaggaagagcgccaggagcccagcgcc 336

acggcccggaaggtggggaggcctggccggaagcgcaagcaccacccgggtggaaagcagt
|||||
acggcccggaaggtggggaggcctggccggaagcgcaagcaccacccgggtggaaagcagt 396

gacacccccaaaggacccagcagtgaccaccaagtctcagcccatggcccaggactctggc
|||||
gacacccccaaaggacccagcagtgaccaccaagtctcagcccatggcccaggactctggc 456

ccctcagatctgctacccaatgggagacttggagaagcggagtgaaccccaacctgaggag
|||||
ccctcagatctgctacccaatgggagacttggagaagcggagtgaaccccaacctgaggag 516

gggagcccagctgcagggcagaaggggtggggccccagctgaaggagaggggaactgagacc
|||||
gggagcccagctgcagggcagaaggggtggggccccagctgaaggagaggggaactgagacc 576

ccaccagaagcctccagagctgtggagaatggctgctgtgtgaccaaggaaggccgtgga
|||||
ccaccagaagcctccagagctgtggagaatggctgctgtgtgaccaaggaaggccgtgga 636

gcctctgcaggagaggggcaaagaacagaagcagaccaacatcgaatccatgaaaatggag
|||||
gcctctgcaggagaggggcaaagaacagaagcagaccaacatcgaatccatgaaaatggag 696

ggctcccggggccgactgcgaggtggcttgggctgggagtccagcctccgtcagcgaccc
|||||
ggctcccggggccgactgcgaggtggcttgggctgggagtccagcctccgtcagcgaccc 756

atgccaagactcaccttccaggcaggggacccctactacatcagcaaacggaaacgggat
|||||
atgccaagactcaccttccaggcaggggacccctactacatcagcaaacggaaacgggat 816

gagtggctggcacgttggaaaaggagggtgagaagaaagccaaggtaatgtagtaatg
|||||
gagtggctggcacgttggaaaaggagggtgagaagaaagccaaggtaatgtagtaatg 876

aatgctgtggaagagaaccaggcctctggagagtctcagaaggtggaggaggccagccct
|||||
aatgctgtggaagagaaccaggcctctggagagtctcagaaggtggaggaggccagccct 936

cctgctgtgcagcagccccacggaccctgcttctccgactgtggccaccacccctgagcca
|||||
cctgctgtgcagcagccccacggaccctgcttctccgactgtggccaccacccctgagcca 996

```

<sup>1</sup> Bolded nucleotides indicate nucleotides that were amended on July 23, 2001.

[illegible]

ctcttgggtggggccaggagctgctcaggcagccattaaggaagacccctggaactgctac  
|||||  
ctcttgggtggggccaggagctgctcaggcagccattaaggaagacccctggaactgctac 1956

atgtgcgggcataagggcacctatgggctgctgcgaagacgggaagactggccttctcga  
|||||  
atgtgcgggcataagggcacctatgggctgctgcgaagacgggaagactggccttctcga 2016

ctccagatgttctttgccaataaccatgaccaggaatttgacccccaaagggtttacca  
|||||  
ctccagatgttctttgccaataaccatgaccaggaatttgacccccaaagggtttacca 2076

cctgtgccagctgagaagaggaagcccatccgcgtgctgtctctctttgatgggattgct  
|||||  
cctgtgccagctgagaagaggaagcccatccgcgtgctgtctctctttgatgggattgct 2136

acagggctcctgggtgctgaaggacctgggcatccaagtggaccgctacattgcctccgag  
|||||  
acagggctcctgggtgctgaaggacctgggcatccaagtggaccgctacattgcctccgag 2196

gtgtgtgaggactccatcacggtgggcatggtgcggcaccagggaaagatcatgtacgtc  
|||||  
gtgtgtgaggactccatcacggtgggcatggtgcggcaccagggaaagatcatgtacgtc 2256

ggggacgtccgcagcgtcacacagaagcatatccaggagtggggcccatcgcacctggtg  
|||||  
ggggacgtccgcagcgtcacacagaagcatatccaggagtggggcccatcgcacctggtg 2316

attggaggcagtccttgcaatgacctctccattgtcaaccctgcccgcaagggactttat  
|||||  
attggaggcagtccttgcaatgacctctccattgtcaaccctgcccgcaagggactttat 2376

gaggggtactggccgcctcttctttgagttctaccgcctcctgcatgatgcgcgcccaag  
|||||  
gaggggtactggccgcctcttctttgagttctaccgcctcctgcatgatgcgcgcccaag 2436

gagggagatgatcgcccccttcttctggctctttgagaatgtggtggccatgggcgttagt  
|||||  
gagggagatgatcgcccccttcttctggctctttgagaatgtggtggccatgggcgttagt 2496

gacaagagggacatctcgcatcttcttgagtctaaccctgcatgattgacgccaaagaa  
|||||  
gacaagagggacatctcgcatcttcttgagtctaaccctgcatgattgacgccaaagaa 2556

gtgtctgctgcacacagggcccggttacttctggggtaaccttcctggcatgaacaggcct  
|||||  
gtgtctgctgcacacagggcccggttacttctggggtaaccttcctggcatgaacaggcct 2616

ttggcatccactgtgaatgataagctggagctgcaagagtgtctggagcacggcagaata  
|||||  
ttggcatccactgtgaatgataagctggagctgcaagagtgtctggagcacggcagaata 2676

gccaagttcagcaaagtgaggaccattaccaccaggtcaaactctataaagcagggcaaa  
|||||  
gccaagttcagcaaagtgaggaccattaccaccaggtcaaactctataaagcagggcaaa 2736

gaccagcatttccccgtcttcatgaacgagaaggaggacatcctgtggtgcaactgaaatg  
|||||  
gaccagcatttccccgtcttcatgaacgagaaggaggacatcctgtggtgcaactgaaatg 2796

gaaaggggtgtttggcttccccgtccactacacagacgtctccaacatgagccgcttggcg  
|||  
gaaaggggtgtttggcttccccgtccactacacagacgtctccaacatgagccgcttggcg 2856

aggcagagactgctgggcccgatcgtggagcgtgccggcatccgccacctcttcgctccg  
|||  
aggcagagactgctgggcccgatcgtggagcgtgccggcatccgccacctcttcgctccg 2916

ctgaaggaatatatttgccttggtgtaagggacatgggggcaaactgaagtagtgatgata  
|||  
ctgaaggaatatatttgccttggtgtaagggacatgggggcaaactgaagtagtgatgata 2976

aaaaagttaaacaacaaacaaacaccaagaacgagaggacggagaaaagtccagcaccc  
|||  
aaaaagttaaacaacaaacaaacaccaagaacgagaggacggagaaaagtccagcaccc 3037

agaagagaaaaaggaatttaaagcaaaccacagaggaggaaaacgccggagggttggcc  
|||  
agaagagaaaaaggaatttaaagcaaaccacagaggaggaaaacgccggagggttggcc 3098

ttgcaaaaggggttgacatcatctcctgagttttcaatgttaaccttcagtcctatctaa  
|||  
ttgcaaaaggggttgacatcatctcctgagttttcaatgttaaccttcagtcctatctaa 3158

aaagcaaaataggccccctccccttcttccccctccggctcctaggaggcgaactttttgttt  
|||  
aaagcaaaataggccccctccccttcttccccctccggctcctaggaggcgaactttttgttt 3218

tctactctttttcagaggggttttctgtttgtttgggtttttgtttcttgctgtgactga  
|||  
tctactctttttcagaggggttttctgtttgtttgggtttttgtttcttgctgtgactga 3278

aacaagagagttattgcagcaaaatcagtaacaacaaaaagtagaaatgccttgagagg  
|||  
aacaagagagttattgcagcaaaatcagtaacaacaaaaagtagaaatgccttgagagg 3338

aaagggagagaggggaaaattctataaaaacttaaaatattgggttttttttttttctt  
|||  
aaagggagagaggggaaaattctataaaaacttaaaatattgggttttttttttttctt 3398

ttctatatatctctttggttgctcttagcctgatcagataggagcacaaacaggaagaga  
|||  
ttctatatatctctttggttgctcttagcctgatcagataggagcacaaacaggaagaga 3458

atagagaccctcggaggcagagtctcctctcccacccccgagcagtcctcaacagcacca  
|||  
atagagaccctcggaggcagagtctcctctcccacccccgagcagtcctcaacagcacca 3518

ttcctgggtcatgcaaaacagaacccaactagcagcagggcgctgagagaaacaccacacca  
|||  
ttcctgggtcatgcaaaacagaacccaactagcagcagggcgctgagagaaacaccacacca 3578

gacactttctacagtatttcaggtgcctaccacacaggaaaccttgaagaaaaccagttt  
|||  
gacactttctacagtatttcaggtgcctaccacacaggaaaccttgaagaaaaccagttt 3638

ctagaagccgctgttacctcttgtttacagtt  
|||  
ctagaagccgctgttacctcttgtttacagtt 3670

# EXHIBIT C

## EXHIBIT C

Alignment spanning the coding region of mouse Dnmt3b from ATCC Deposit No. 209934 (top) and currently amended SEQ ID NO:2 (bottom)<sup>2</sup>

```

caggaaacaatgaagggagacagcagacatctgaatgaagaagaggggtgccagcgggtat
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
caggaaacaatgaagggagacagcagacatctgaatgaagaagaggggtgccagcgggtat 319

gaggagtgcattatcgттаатgggaacttcagtгaccagtcctcagacacgaaggatgct
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
gaggagtgcattatcgттаатgggaacttcagtгaccagtcctcagacacgaaggatgct 379

ccctcacccccagtccttgaggcaatctgcacagagccagtcctgcacaccagagaccaga
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
ccctcacccccagtccttgaggcaatctgcacagagccagtcctgcacaccagagaccaga 439

ggccgcaggtcaagctcccggctgtctaagagggaggtctccagccttctgaattacacg
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
ggccgcaggtcaagctcccggctgtctaagagggaggtctccagccttctgaattacacg 499

caggacatgacaggagatggagacagagatgatgaagtagatgatgggaatggctctgat
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
caggacatgacaggagatggagacagagatgatgaagtagatgatgggaatggctctgat 559

attctaatgccaaagctcacccgtgagaccaaggacaccaggacgcgctctgaaagcccg
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
attctaatgccaaagctcacccgtgagaccaaggacaccaggacgcgctctgaaagcccg 619

gctgtccgaacccgacatagcaatgggacctccagcttgгagaggcaaagagcctcccc
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
gctgtccgaacccgacatagcaatgggacctccagcttgгagaggcaaagagcctcccc 679

agaatcacccgaggtcggcagggccgcccacatgtgcaggagtaccctgtggagtttccg
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
agaatcacccgaggtcggcagggccgcccacatgtgcaggagtaccctgtggagtttccg 739

gctaccaggtctcgggacgtcgagcatcgтcttcagcaagcacgccatggтcatccct
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
gctaccaggtctcgggacgtcgagcatcgтcttcagcaagcacgccatggтcatccct 799

gccagcgtcgacttcatggaagaagtгacacctaagagcgtcagtagcccatcagttgac
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
gccagcgtcgacttcatggaagaagtгacacctaagagcgtcagtagcccatcagttgac 859

ttgagccaggatggagatcaggaggggtatggataccacacaggtggatgcagagagcaga
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
ttgagccaggatggagatcaggaggggtatggataccacacaggtggatgcagagagcaga 919

gatggagacagcacagagtatcaggatgataaagagtttggaataggtgacctcgтgtgg
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
gatggagacagcacagagtatcaggatgataaagagtttggaataggtgacctcgтgtgg 979

ggaaagatcaagggcttctcctggтggcctgccatggтggтgtcctggaaagccacctcc
|||||||||||||||||||||||||||||||||||||||||||||||||||||||||
ggaaagatcaagggcttctcctggтggcctgccatggтggтgtcctggaaagccacctcc 1039

```

<sup>2</sup> Bolded nucleotides indicate nucleotides that were amended on July 23, 2001.

aagcgacaggccatgcccggaatgcgctgggtacagtgggttggtgatggcaagttttct  
|||  
aagcgacaggccatgcccggaatgcgctgggtacagtgggttggtgatggcaagttttct 1099

gagatctctgctgacaaactgggtggctctggggctgttcagccagcactttaatctggct  
|||  
gagatctctgctgacaaactgggtggctctggggctgttcagccagcactttaatctggct 1159

accttcaataagctggtttcttataggaaggccatgtaccacactctggagaaagccagg  
|||  
accttcaataagctggtttcttataggaaggccatgtaccacactctggagaaagccagg 1219

gttcgagctggcaagaccttctccagcagtcctggagagtcactggaggaccagctgaag  
|||  
gttcgagctggcaagaccttctccagcagtcctggagagtcactggaggaccagctgaag 1279

cccatgctggagtgggcccacggtggcttcaagcctactgggatcgagggcctcaaacc  
|||  
cccatgctggagtgggcccacggtggcttcaagcctactgggatcgagggcctcaaacc 1339

aacaagaagcaaccagtggttaataagtccaaggtgcgtcgttcagacagtaggaactta  
|||  
aacaagaagcaaccagtggttaataagtccaaggtgcgtcgttcagacagtaggaactta 1399

gaacccaggagacgcgagaacaaaagtccaagacgcacaaccaatgactctgctgcttct  
|||  
gaacccaggagacgcgagaacaaaagtccaagacgcacaaccaatgactctgctgcttct 1459

gagtccccccacccaagcgcctcaagacaaatagctatggcggaaggaccgaggggag  
|||  
gagtccccccacccaagcgcctcaagacaaatagctatggcggaaggaccgaggggag 1519

gatgaggagagccgagaacggatggcttctgaagtcaccaacaacaagggaatctggaa  
|||  
gatgaggagagccgagaacggatggcttctgaagtcaccaacaacaagggaatctggaa 1579

gaccgctgtttgtcctgtggaagaagaaccctgtgtccttccacccctctttgaggg  
|||  
gaccgctgtttgtcctgtggaagaagaaccctgtgtccttccacccctctttgaggg 1639

gggctctgtcagagttgccgggatcgcttcttagagctcttctacatgtatgatgaggac  
|||  
gggctctgtcagagttgccgggatcgcttcttagagctcttctacatgtatgatgaggac 1699

ggctatcagtcctactgcaccgtgtgctgtgagggccgtgaactgctgctgtgcagtaac  
|||  
ggctatcagtcctactgcaccgtgtgctgtgagggccgtgaactgctgctgtgcagtaac 1759

acaagctgctgcagatgcttctgtgtggagtgtctggaggtgctggtgggcgcaggcaca  
|||  
acaagctgctgcagatgcttctgtgtggagtgtctggaggtgctggtgggcgcaggcaca 1819

gctgaggatgccaaagctgcaggaaccctggagctgctatatgtgcctccctcagcgctgc  
|||  
gctgaggatgccaaagctgcaggaaccctggagctgctatatgtgcctccctcagcgctgc 1879

catggggctcctccgacgcaggaaagattggaacatgcgcctgcaagacttcttcactact  
|||  
catggggctcctccgacgcaggaaagattggaacatgcgcctgcaagacttcttcactact 1939



```

gatcctgacctggaagaatttgagccaccaagttgtacccagcaattcctgcagccaaa
|||||
gatcctgacctggaagaatttgagccaccaagttgtacccagcaattcctgcagccaaa 1999

aggaggcccattagagtcctgtctctgtttgatggaattgcaacgggggtacttggtgctc
|||||
aggaggcccattagagtcctgtctctgtttgatggaattgcaacgggggtacttggtgctc 2059

aaggagttgggtattaaagtggaaaagtacattgcctccgaagtctgtgcagagtcctac
|||||
aaggagttgggtattaaagtggaaaagtacattgcctccgaagtctgtgcagagtcctac 2119

gctgtgggaactgttaagcatgaaggccagatcaaatatgtcaatgacgtccggaaaatc
|||||
gctgtgggaactgttaagcatgaaggccagatcaaatatgtcaatgacgtccggaaaatc 2179

accaagaaaaatattgaagagtggggcccgttcgacttggtgattggtggaagcccatgc
|||||
accaagaaaaatattgaagagtggggcccgttcgacttggtgattggtggaagcccatgc 2239

aatgatctctctaacgtcaatcctgcccgcgaagggtttatatgagggcacaggaaggctc
|||||
aatgatctctctaacgtcaatcctgcccgcgaagggtttatatgagggcacaggaaggctc 2299

ttcttcgagttttaccacttgctgaattatacccgccccaaggagggcgacaaccgtcca
|||||
ttcttcgagttttaccacttgctgaattatacccgccccaaggagggcgacaaccgtcca 2359

ttcttctggatgttcgagaatgttggtggccatgaaagtgaatgacaagaaagacatctca
|||||
ttcttctggatgttcgagaatgttggtggccatgaaagtgaatgacaagaaagacatctca 2419

agattcctggcatgtaaccagtgatgatcgatgccatcaagggtgtctgctgctcacagg
|||||
agattcctggcatgtaaccagtgatgatcgatgccatcaagggtgtctgctgctcacagg 2479

gcccgggtacttctggggtaacctaccgggaatgaacaggcccgtgatggcttcaaagaat
|||||
gcccgggtacttctggggtaacctaccgggaatgaacaggcccgtgatggcttcaaagaat 2539

gataagctcgagctgcaggactgcctggagttcagtaggacagcaaagttaaagaaagtg
|||||
gataagctcgagctgcaggactgcctggagttcagtaggacagcaaagttaaagaaagtg 2599

cagacaataaccaccaagtgcgaactccatcagacagggcaaaaaccagcttttccctgta
|||||
cagacaataaccaccaagtgcgaactccatcagacagggcaaaaaccagcttttccctgta 2659

gtcatgaatggcaaggacgacgttttgtggtgcactgagctcgaaaggatcttcggcttc
|||||
gtcatgaatggcaaggacgacgttttgtggtgcactgagctcgaaaggatcttcggcttc 2719

cctgctcactacacggacgtgtccaacatggggccgcgcccgtcagaagctgctgggc
|||||
cctgctcactacacggacgtgtccaacatggggccgcgcccgtcagaagctgctgggc 2779

aggtcctggagtgtaccggtcatcagacacctgtttgcccccttgaaggactactttgcc
|||||
aggtcctggagtgtaccggtcatcagacacctgtttgcccccttgaaggactactttgcc 2839

```

tgtgaatagttctacccaggactggggagctctcggtcagagccagtgtcccagagtc  
||||||||||||||||||||||||||||||||||||||||||||||||||||||  
tgtgaatagttctacccaggactggggagctctcggtcagagccagtgtcccagagtc 2896

# EXHIBIT D

LETTER

# Estimation of Errors in "Raw" DNA Sequences: A Validation Study

Peter Richterich<sup>1</sup>

Genome Therapeutics Corp., Waltham, Massachusetts 02154 USA

As DNA sequencing is performed more and more in a mass-production-like manner, efficient quality control measures become increasingly important for process control, but so also does the ability to compare different methods and projects. One of the fundamental quality measures in sequencing projects is the position-specific error probability at all bases in each individual sequence. Accurate prediction of base-specific error rates from "raw" sequence data would allow immediate quality control as well as benchmarking different methods and projects while avoiding the inefficiencies and time delays associated with resequencing and assessments after "finishing" a sequence. The program PHRED provides base-specific quality scores that are logarithmically related to error probabilities. This study assessed the accuracy of PHRED's error-rate prediction by analyzing sequencing projects from six different large-scale sequencing laboratories. All projects used four-color fluorescent sequencing, but the sequencing methods used varied widely between the different projects. The results indicate that the error-rate predictions such as those given by PHRED can be highly accurate for a large variety of different sequencing methods as well as over a wide range of sequence quality.

In DNA sequencing, knowledge about the accuracy of sequences can be very valuable. For example, different large-scale sequencing projects may produce sequences at similar rates and costs but with significantly different error rates in the final sequence. One major determinant in the final error rate is the accuracy of the "raw" sequence. Knowledge about the frequency and location of errors in the raw sequence data can help to direct "polishing" efforts to the places where additional effort is needed; it also enables the comparison between different sequencing projects without requiring that the same region be sequenced in each project.

Another area where estimates about sequence error rates would be beneficial is technology development. Accurate error estimates at each base would enable "quality benchmarking" between different methods, thus enabling researchers to choose the method that fills their needs for accuracy and throughput best.

Several groups have developed mathematical models to predict the error probability at any given position in raw sequences. Lawrence and Solovyev used linear discriminant analysis to calculate separate probability estimates for insertions, deletions, and mismatches (Lawrence and Solovyev 1994). Iwring and Green (1998) developed the program

PHRED, which calculates a quality score at each base. This quality score  $q$  is logarithmically linked to the error probability  $p$ :  $q = -10 \times \log_{10}(p)$  (for a discussion of how quality scores are calculated and what the limitations are, see Ewing et al. (1998). When used in combination with sequence assembly and finishing programs that utilize these error estimates, reliable error probabilities promise to increase the accuracy of consensus sequences and to reduce the efforts required in the finishing phase of sequencing projects (Churchill and Waterman 1992; Bonfield and Staden 1995).

To examine the accuracy of probability estimates made by the program PHRED, we compared the actual and predicted error rates for six different cosmid- or BAC-sized projects that were produced by six different large-scale sequencing centers in the United States. All of these six projects used four-color fluorescent sequencing machines; however, the DNA preparation methods, sequencing enzymes, fluorescent dyes and chemistries, and gel lengths varied significantly between the six groups. Table 1 gives an overview of the sequencing projects analyzed. Table 2 lists the different methods used.

## RESULTS

### Error Rate Prediction Accuracy for Six Projects

A comparison of actual and predicted error rates for the six projects in this study is shown in Table 3.

<sup>1</sup>E-MAIL [peter.richterich@genomecorp.com](mailto:peter.richterich@genomecorp.com); FAX (781) 893-9535.

**Table 1. Summary of Data Sets**

Project	Reads	Aligned bases	Average aligned read length
A	455	416,214	915
B	1277	871,230	682
C	1065	603,655	567
D	834	414,595	497
E	1638	1,149,209	702
F	1885	907,796	482
Total	7154	4,362,699	610

The results indicate that PHRED is very successful in identifying bases with low error probabilities. For example, the 1.28 million bases with quality scores of 4–12 (corresponding to error probabilities between 39.8% and 6.3%) contain a total of 187,926 errors. In contrast, the 1.44 million bases with quality scores between 33 and 42 (corresponding to error probabilities between 0.05% and 0.006%) contain only 237 errors, which translates into a 790-fold lower error rate. The trend toward lower error rates can also be observed for each individual project. In most cases, the actual number of errors is close to the predicted error rate. It is also apparent that the actual error rate is typically lower than the predicted error rate.

Both the high overall accuracy and the tendency to slightly overpredict errors are confirmed by statistical analysis, as shown in Table 4. The correlation between predicted and actual error frequencies is excellent for all projects (Spearman correlation coefficient  $>0.89$ ,  $P < 0.0001$ ). Averaged over all projects, the actual error rate is 84.5% of the predicted error rate; the slope of the relation between predicted and actual error rates differs slightly between projects and ranges from 76.6% to 88.4%. To put these differences between projects in relation, it is worthwhile remembering that PHRED quality scores cover a wide dynamic range: The maximum quality score of 51 corresponds to a 50,000-fold lower predicted error rate than the minimum quality score of 4. Even the relative difference between successive quality is larger than the relative difference in the slopes; for example, a quality score of 10 corresponds to an error probability of 10%, whereas a score of 9 corresponds to an error probability of 12.6%.

A different way of looking at the relation between the actual and predicted error rates is shown

in Figure 1. Here, the error rates as a function of the position within all reads in each of the projects, averaged over 50-base windows, is depicted. For all six projects, the predicted error rates are very close to the actual error rates over the entire length of the sequences. Each project has a characteristic distribution of error rates, which differs from each of the other projects. The minimum error rate differs dramatically between projects. The best projects achieve raw error rates of 0.23%–0.36% in the best region of the sequence read, typically from base 150 to 200. The worst project in the data set had an ~10-fold higher error rate of 2.58%.

Toward the end of sequence reads, the error rates increase and start to exceed 10% between bases 300 and 700. In projects that used mainly short gels (e.g., projects D and F), this increase begins sooner, whereas projects that use longer gels show a markedly longer stretch of low error rates (e.g., projects A and B).

Table 5 summarizes key results for the six projects. The first four projects have similar minimum and average error rates. However, the length of the region where the error rate is below 5% differs significantly, from 403 to 682 bases. The project with the shorter low error rate regions contained larger portions of reads generated on short gels, whereas projects A and B were run exclusively on long gels (ABI373 stretch or ABI377 sequencers). Other factors contributing to differences between the first four projects were differences in sequencing chemistries, production scale, and electrophoresis conditions and machines.

Project E and, in particular, project F, had significantly higher error rates than the first four projects. In projects E and F, every sequence generated for the project had been included in the data set, whereas the other four projects had eliminated some "bad" sequences through manual or auto-

**Table 2. Overview of Sequencing Methods Used in the Different Projects**

Template DNA	single-stranded M13; double-stranded plasmids
Sequencing enzymes	Sequenase, <i>Taq</i> , KlenTaqTR, AmpliTaq FS
Sequencing chemistries	Dyes primer (two different dyes chemistries), dye terminator
Sequencing machines	ABI 373, ABI 373 stretch, ABI 377
Gel length	Only short gels, only long gels, mixes of short and long gels

**Table 3. Comparison of Predicted and Actual Error Rates for Six Different Sequencing Projects**

Project	Quality score	4-12	13-22	23-32	33-42	43-51
A	aligned bases	119,246	75,293	70,391	144,876	73,234
	expected errors	20,256	2,064	172	37	1
	actual errors	16,784	1,758	127	17	1
B	aligned bases	182,034	137,940	181,998	399,690	140,176
	expected errors	29,953	3,704	410	102	3
	actual errors	26,038	2,536	287	35	0
C	aligned bases	139,345	131,419	151,197	292,070	68,529
	expected errors	22,277	3,411	357	74	2
	actual errors	16,670	1,513	194	26	3
D	aligned bases	103,898	68,995	68,613	153,730	111,752
	expected errors	16,880	1,919	168	38	3
	actual errors	14,495	1,924	146	59	2
E	aligned bases	378,755	217,438	167,968	392,717	144,313
	expected errors	63,947	6,336	418	95	4
	actual errors	55,968	6,516	355	67	5
F	aligned bases	359,809	136,688	98,840	64,035	5,130
	expected errors	66,938	4,079	256	23	0
	actual errors	57,971	3,856	332	33	1
All	aligned bases	1,283,087	767,773	739,007	1,447,118	543,134
	expected errors	220,252	21,513	1,781	370	13
	actual errors	187,926	18,103	1,441	237	12

matic inspection. After eliminating <10% of the worst sequences in project E, the error rates for the remaining sequences were comparable to those of the first four projects. In contrast, project F showed a much more uniform distribution of sequence quality.

The last column in Table 5 shows the average number of bases with an estimated error probability of at most 0.1%, which is equivalent to a quality score of at least 30. The count of such "very high-quality" bases is a good indicator of sequence quality, both for individual sequences and, when aver-

**Table 4. Summary of Statistical Analysis Results**

Project	Spearman $\rho$	$P >  \rho $	Slope	t ratio	$P >  t $
A	0.9646	<0.0001	0.818	75.1	<0.0001
B	0.9890	<0.0001	0.874	98.2	<0.0001
C	0.9846	<0.0001	0.766	71.6	<0.0001
D <sup>a</sup>	0.8692	<0.0001	0.855	68.3	<0.0001
E	0.9956	<0.0001	0.884	144.3	<0.0001
F	0.9968	<0.0001	0.865	151.6	<0.0001
All	0.9964	<0.0001	0.845	174.5	<0.0001

<sup>a</sup>In project D, the Spearman correlation coefficient  $\rho$  was artificially low as only very few bases (10) bases had a quality score of 5, and none of these bases contained an actual error (expected: 3.16 errors). Exclusion of this quality score gave a Spearman correlation coefficient of 0.9786 ( $P < 0.0001$ ). The frequencies in the slope calculations were weighed by the number of bases at any given quality score and, thus, were not sensitive to such small sample distortions (see Methods).

ity analysts and control in large-scale DNA sequencing projects. To analyze how accurate PHRED error estimates are for different quality sequences within the same sequencing project, we subdivided a data set into four quartiles, based on the number of very high-quality bases in each sequence (see Methods). The comparison of actual and predicted error rates is shown in Figure 2.

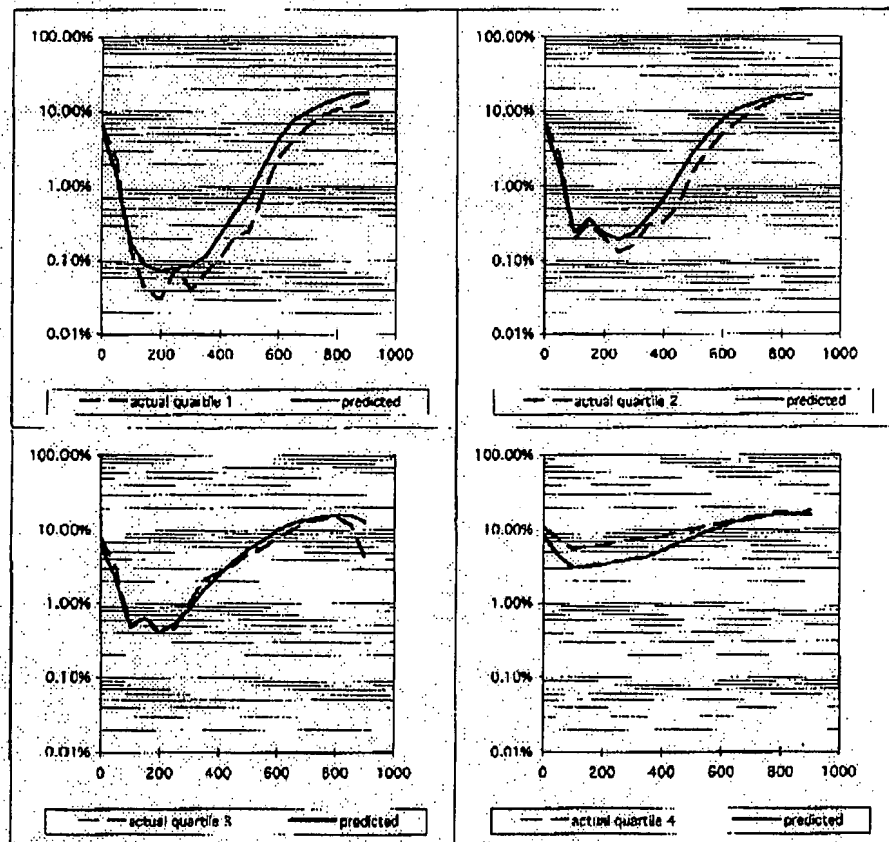
When measured by the error rate in the best region of a sequence, the data quality in the different quartiles varies >100-fold between the best and the worst 25% of the sequences. The best quartile showed ~0.03% error for >100 bases, whereas the error rate in the worst quartile always exceeded 5%. In quartiles 2 and 3, the predicted error rates match the actual error rates very closely. In the best and

worst quartiles, PHRED's accuracy was somewhat lower from base 100 to 500. In the best sequences, PHRED's error estimates were about twofold too high; in the worst sequences, the error estimates were too low, again by a factor of 2. This underprediction of errors can be partially explained by the fact that PHRED gives ambiguous base calls (N's) a quality score of 4, corresponding to an error probability of 39.8%; however, N's will always show up as an actual error. Even in the worst and best quartiles, however, the predicted error rate curves are very similar to the actual error rate curves.

The results shown in Figure 2 also demonstrate that the count of very high-quality bases, or bases with an estimated error probability of at most 0.1%, can be used effectively to characterize the overall

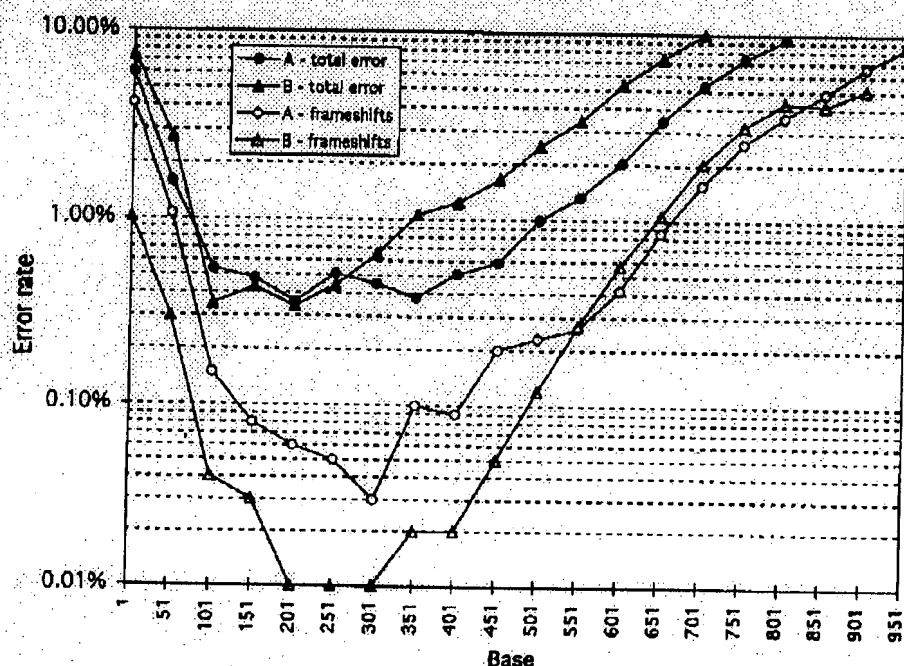
quality of a sequence read. Sorting the sequence reads into quartiles based on the number of very high-quality bases worked well, as shown by the >100-fold difference in the minimum error rate between the first and the fourth quartile.

Other methods to characterize the overall quality of individual reads based on PHRED quality scores can give similar results. For example, counting bases above a minimum quality threshold anywhere in the range of 20–40 gave similar results for most data sets (not shown), and such counts are used by a number of different laboratories as quality measures. Alternatively, the quality values can be converted to error probabilities and averaged to give the predicted error rate for the trace, or summed to give the total predicted number of errors in a trace. However, such averages and totals can sometimes give a misleading picture, as the following example illustrates. Assume that two sequence reads have very similar quality in the alignable part of the read but that one of the two sequences was run much longer and



**Figure 2** Actual and predicted error rates in different quality subsets of project B. Sequence reads were sorted by the number of bases with a predicted error rate of at most 0.1% (very high-quality bases), and assigned to quartiles, with quartile 1 corresponding to the highest numbers. Actual and predicted error rates for all sequences in each subset were calculated as in Fig. 1. Note that a number of sequence reads that had been rejected because of too low quality were added back to the data set for illustrative purposes, all of which are in quartile 4. These sequences were not included in the data sets used to generate Figs. 1 and 3 and Tables 1 and 3.

## RICHTERICH



**Figure 3** Actual frameshift and total error rates for projects A and B. To calculate frameshift error rates, only insertions and deletions were counted. Mismatch errors, which account for the vast majority of errors after base 150, were included only in the total error count. Note that project B ( $\Delta$ ,  $\triangle$ ) has a slightly similar or slightly higher total error rate compared to project A ( $\bullet$ ,  $\circ$ ) but only about one-third as many insertions and deletions up to base 500. For both projects, the frameshift error rate in the raw data is  $<1$  in 1000 for  $>300$  bases, and  $\leq 1$  in 10,000 for  $>100$  bases in project B.

therefore contains a longer unalignable "tail" of very low-quality bases. When calculating the average error rate for these two sequences, the second sequence will have a much higher average error and, therefore, appear to be of lower quality. In contrast, the counts of very high-quality bases for both sequences will be very similar, as the unalignable tails contain few, if any, high-quality bases. Therefore, counts of bases above a high enough quality threshold will give a more robust and clearer picture of trace quality.

### Frameshift Error Rates for Different Sequencing Chemistries

Depending on how biologists use DNA sequences, knowledge about total error rates in raw sequences may or may not be sufficient. For example, frameshift errors in coding sequences will generally lead to incorrectly predicted open reading frame, whereas mismatch errors will do so only if the mismatch introduces a stop codon or a new splice site. At the time of this writing, PHRED did not differentiate between mismatch and frameshift errors, but only estimated total error rates. This might occa-

sionally lead to questionable conclusions, as the results shown in Figure 3 illustrate.

Figure 3 shows the total actual error rates and the frameshift error rates for two projects, A and B. The total error rates for both projects are similar for up to 350 bases; after 350 bases, project B has a somewhat higher total error rate. However, examining the frameshift error rate gives rise to a different picture: from base 1 to 500, project A has approximately four times as many insertions and deletions as project B. This difference in frameshift error rates can be explained by the sequencing chemistries that were used in the two projects. Project B, with the lower frameshift error rate, used only dye terminator chemistry, which is known to eliminate band spacing artifacts from hairpin structures ("compressions"). Project A, on the

other hand, used dye primer chemistry, which is more prone to insertion and deletion errors from mobility artifacts, for most sequencing reactions.

## DISCUSSION

As large-scale DNA sequencing has become a more routine and common process, the traditional methods for assessing sequence quality have become unsatisfactory. In projects like single-pass cDNA sequencing, it is not possible to calculate and compare error rates after finishing a sequence, as finishing never takes place. Even when a comparison between raw and finished sequence can be done, the time delay between raw data generation and quality assessment is often large. This delay makes it difficult to improve ongoing projects, and it sometimes makes it impossible to capture problems early on. Some immediate quality feedback can be reached by including known standard sequences for quality control. However, this approach can be costly, and it fails when error profiles differ between standard and unknown sequences.

In contrast to these traditional methods to assess sequence accuracy, direct estimation of error



rates in raw sequence data would enable immediate quality control and feedback. Accurate, base-by-base estimates of error probabilities could also increase the utility of single-pass sequences significantly, allow efficient comparison and optimization of different sequence chemistries, and enable the development of better software tools for sequence assembly and analysis.

The critical question for any error rate prediction tool is how accurate are the error rate estimates, in particular if different sequencing methods and chemistries are used? The results presented herein provide an answer to this question for the program PHRED, as well as clues where further development would be useful. As shown in Tables 3 and 4 and in Figure 1, the agreement between predicted and actual error rates was very good in each of the six different projects analyzed. The observed high level of prediction accuracy in all of these projects is almost astonishing if one takes into account that actual errors are binary (a base is either correct or wrong), whereas predicted error rates are probabilities on a scale from 0.0 to 1.0. The observed tendency to overpredict error rates can be at least partially explained by the "small sample correction" that was used in the derivation of threshold parameters for quality scores (Iwling and Green 1998). For most practical applications, such a somewhat conservative estimation of quality scores is tolerable or even desirable. Overall, the results clearly show that error probabilities given by PHRED accurately describe raw sequence data quality.

In judging the usefulness of predicted error probabilities, it is important to know how differences in sequencing methods will influence the prediction accuracy. For example, the larger variation in peak heights tends to be larger in dye terminator sequencing than in dye primer sequencing, and different sequencing enzymes are known to produce different specific height variation patterns. Any estimation of error probabilities that takes the peculiarities of a specific sequencing chemistry into account would therefore be expected to be less accurate for different chemistries.

The projects included in this study were specifically chosen to provide an initial answer to the question of how generally useful PHRED quality scores are. These projects represent the vast majority of different multicolor fluorescent sequencing methods used in the last 3 years: different template DNAs and DNA preparation methods, different enzymes, gel lengths, run conditions, and different fluorescent dyes. The data also include a considerable spread in data quality, both between projects

and within individual projects. None of the projects analyzed here were included in PHRED's training set, and just one of the six laboratories that contributed data to this study also contributed data to the training data sets. One of the projects in this study consisted entirely of dye terminator sequences, which presented only a small fraction of the sequences in the test data set. Another project exclusively used a set of fluorescent dyes different from those used in the training sets. Each project differed from the other projects in this study in at least one, and typically many, experimental aspects like template preparation, sequencing enzymes, gel run conditions, and so forth. Despite these differences, the accuracy of error rate predictions was very similar for all projects.

Our results justify some optimism about the accuracy of PHRED quality scores for minor changes in sequencing technology, for example, sequences generated by new enzymes and fluorescent dyes. Initial studies showed that PHRED quality scores were also accurate for sequences produced by multiplex sequencing with radioactive detection (P. Richterich, unpubl.). However, we also observed two effects that can invalidate PHRED quality scores during these studies. First, sequences generated by chemical sequencing gave too low quality scores at mixed (A + G) reactions. Because secondary peak height is one of the parameters used in the error rate predictions, this is not surprising. Another potential source of error is high-frequency noise in the trace data. With such data, PHRED occasionally underestimated the band spacing by a factor of 2 or more, which resulted in incorrect base calls and quality scores. By applying simple smoothing algorithms to data with high-frequency noise, these problems could typically be resolved. Similar steps may be necessary to obtain accurate PHRED quality scores on data that have been generated by different sequencing instruments or preprocessed by different software.

Accurate quality scores can have a major impact on how sequences are used downstream from the sequence production process. In traditional sequencing projects where the goal is complete coverage at a final error rate below (e.g.) 1 in 10,000, the accuracy goals can be reached with single sequence reads as long as the quality scores are at least 40 (however, other potential problems like clone instability may make higher coverage advisable). Interesting questions arise as to how individual read quality contributes to project quality, or the error rate of the "final" sequence. Under the assumption that errors between different sequence reads are

## RICHTERICH

completely independent, one could argue that two reads with a quality score of 20 (error probability of 1 in 100) are just as valuable as one sequence with a quality score of 40 (error probability of 1 in 10,000). However, although a single sequence stretch with quality levels above 40 would give a final sequence with an error rate of <1 in 10,000, assembling a consensus from two sequences with quality scores of 20 (1% error rate) could lead to one of two results: If the errors were completely random, the consensus sequence would be ambiguous at 2% of all locations; if the errors were completely localized, for example, because of reproducible compressions, the consensus sequence would have one "hidden" error every 100 bases. Typically, consensus sequences derived from low-quality sequences will have both kinds of problematic regions. Increased coverage can rapidly eliminate the random errors; however, increased coverage does not resolve errors from systematic sources. Manual examination of such problem areas is generally required; such "contig editing," however, tends to be time consuming, requires highly trained personnel, is an obstacle toward complete automation of DNA sequencing, and sometimes fails to eliminate all errors. This leads to the somewhat counterintuitive conclusion that the practical value of increasing sequence quality can be even higher than indicated by the quality scores. One sequence of average quality above 40 can be "worth" more than two sequences of average quality 20.

Another application of DNA sequencing where high quality can be of disproportionately high value is the search for mutations in genomic DNA. In low quality sequences, secondary peaks and low resolution often complicate the identification of heterozygous mutations. In regions of higher sequence quality, such secondary peaks are smaller or absent and peaks are better resolved. Therefore, both false-positive and false-negative errors can be significantly reduced in high-quality regions. Tools like PHRED, which can accurately measure sequence quality from trace data, can be of twofold value for mutation detection. First, base-specific quality scores can allow optimization of sequencing methods and strategies for mutation detection. Second, the quality scores can be used to evaluate the usefulness of individual sequence reads for mutation detection (e.g., by discarding reads below minimum thresholds), and they can guide software that automatically detects mutations.

The ability to predict error rates in a highly accurate fashion is likely to have a major impact in applications like those described above. PHRED is

the first widely used program that accurately predicts base-specific error probabilities. However, the algorithm for determining quality values has been described (Ewing and Green 1998), and it should be straightforward to implement similar quality values in other base-calling programs. Furthermore, an extension of the approach developed by Ewing and Green should be possible. For example, differentiation between mismatch and frameshift errors would enable better comparisons of sequencing methods with similar total error rates but different frameshift error rates. Several groups have described efforts to calculate separate probabilities (or "confidence assessments") for mismatch errors and frameshift errors (Lawrence and Solov'yev 1994; Berno 1996). Their results demonstrated that different approaches to error type characterization are feasible and promising. Implementation of such error type predictions in other programs similar to the way PHRED uses quality scores would enable better method assessments, benchmarking, and production quality control, and could have a significant impact on downstream uses of DNA sequence information.

## METHODS

### Data Sets

For one project, sequence raw data in the form of ABI trace files were downloaded from a public FTP site. Sequence data for the five other projects were kindly provided by five different large-scale sequencing groups. Table 1 gives a summary of the six projects, and Table 2 gives an overview of the different sequencing methods used in the projects. The projects differed in the amount of prescreening of data that had been done, reflecting different approaches to quality control in different laboratories. In two projects (B and C), different software programs had been used to identify and eliminate low-quality sequences. One project (F) included all data files generated, whereas the other three projects had excluded "failed lanes."

### Comparison of Actual and Predicted Error Rates

The sequences for all traces in each project were recalled using the program PHRED (v. 961028). Next, sequences in each project were assembled with PHRAP (P. Green, unpubl.). Slightly different methods were chosen for the statistical and graphical evaluation of the error rate prediction accuracy. In the statistical evaluation, only the longest contig produced by PHRAP was considered. The tables of aligned bases and observed discrepancy counts for

each quality score were taken from the PHRAP output and analyzed as follows. The expected number of discrepancies ( $E$ ) at each quality score ( $q$ ) was calculated by multiplying the number of aligned bases ( $N$ ) with the error probability corresponding to the quality score:  $E = N 10^{-0.1q}$ . The Spearman ranking coefficients were calculated by comparing the expected and observed error frequencies. To obtain the quantitative relation between the expected and observed error rates over the entire range, a least-squares fit between the observed and expected rates was performed, with the intercept set to zero and the number of aligned bases at each quality score used as weights.

For a graphical comparison of estimated and actual error rates in 50-bp windows, the following steps were taken. For two of the projects, the consensus sequence was retrieved from public databases. For the four other projects, the DNA sequence and quality information were used by the program PHRAP to assemble consensus sequences for each of the projects. The individual reads were aligned to the consensus sequences of the longest contig, using the program CROSS\_MATCH (P. Green, unpubl.), after removing single-coverage regions from the ends of the consensus sequence. CROSS\_MATCH uses an implementation of the Smith-Waterman algorithm to generate alignments that typically do not include the ends of sequences, where disagreements are commonly due to vector sequence or low quality sequence.

The quality files generated by PHRED and the alignment summaries generated by CROSS\_MATCH were then analyzed as follows. First, the region of each query sequence that had been aligned by CROSS\_MATCH was determined. Next, the actual and predicted error rates for the entire aligned part of each individual sequence was calculated. In addition, the average actual and predicted error rates for all alignable sequences together were calculated for windows of 50 bases in length. To calculate the predicted error rate, the quality scores  $q$  determined by PHRED at each base were converted to error probabilities as described above (Ewing and Green 1998).

#### Subdividing Data Into Subsets Based on Data Quality

To examine the accuracy of PHRED quality scores for data subsets of different quality within a project, the following approach was taken. For all sequence reads in project B, the number of bases with a quality score of at least 30 in each sequence was determined (bases with quality scores of at least 30 were called very high-quality bases, or VHQ bases). Se-

quences were sorted in descending order based on the number of very high-quality bases, and divided into four quartiles. Accordingly, quartile 1 contained 25% of sequences with the highest number of very high-quality bases, and quartile 4 contained the "worst" sequences. To illustrate the prediction accuracy in data with relatively high error rates, sequences from project B that had been "discarded" because they had not met the minimum quality criteria were added back to the data set. The sequences in each quartile were compared to the consensus sequences that had been generated using the entire data set, as described above for the graphical comparison.

#### Determining Actual Frameshift Error Rates

The calculation of actual frameshift error rates in the raw sequence data was performed using CROSS\_MATCH, similar to the procedure described above for total error rates, except that only insertion and deletion errors were counted. Because PHRED does not give separate frameshift error estimates, a comparison of predicted and actual frameshift errors is not possible.

#### ACKNOWLEDGMENTS

I thank the participating laboratories for contributing their data; Dr. Josée Dupuis for help with the statistical analysis, and Dr. Phil Green for helpful discussions.

The publication costs of this article were defrayed in part by payment of page charges. This article must therefore be hereby marked "advertisement" in accordance with 18 USC section 1734 solely to indicate this fact.

#### REFERENCES

- Berno, A.J. 1996. A graph theoretic approach to the analysis of DNA sequencing data. *Genome Res.* 6: 80-91.
- Bonfield, J.K. and R. Staden. 1995. The application of numerical estimates of base calling accuracy to DNA sequencing projects. *Nucleic Acids Res.* 23: 1406-1410.
- Churchill, G. and M.S. Waterman. 1992. The accuracy of DNA sequences: estimating sequence quality. *Genomics* 14: 89-98.
- Ewing, R. and P. Green. 1998. Base-calling of automated sequencer traces using phred. II. Error probabilities. *Genome Res.* (this issue).
- Ewing, B., L. Hillier, M.C. Wendt, and P. Green. 1998. Base-calling of automated sequencer traces using phred. I. Accuracy assessment. *Genome Res.* (this issue).
- Lawrence, C.B. and V.V. Solovyev. 1994. Assignment of position-specific error probability to primary sequence data. *Nucleic Acids Res.* 22: 1272-1280.

Received October 27, 1997; accepted in revised form February 3, 1998.

# GENOME RESEARCH

Volume 8 Number 3  
March 1998

## Editors

**Laurie Goodman**  
Cold Spring Harbor Laboratory  
**Mark Boguski**  
National Center for Biotechnology  
Information, NIH  
**Aravinda Chakravarti**  
Case Western Reserve University

## Editorial Board

**Rakesh Anand**  
Zeneca Pharmaceuticals  
**Styllanos Antonarakis**  
University of Geneva  
**Charles Auffray**  
CNRS  
**Philip Avner**  
Institut Pasteur  
**Andrew Ballabio**  
Telethon Institute of Genetics and  
Medicine  
**David Bentley**  
The Sanger Centre  
**Bruce Birren**  
Whitehead Institute/MIT Center for  
Genome Research  
**Michael Boehnke**  
University of Michigan School of  
Public Health  
**Anne Bowcock**  
University of Texas Southwestern  
Medical Center  
**David Burke**  
University of Michigan Medical School  
**Jeffrey Chamberlain**  
University of Michigan Medical School  
**Elson Chen**  
Perkin-Elmer Corporation  
**David R. Cox**  
Stanford University School of Medicine  
**Ronald W. Davis**  
Stanford University School of Medicine  
**Richard Durbin**  
Sanger Centre, UK  
**Joseph Ecker**  
University of Pennsylvania  
**Beverly S. Emanuel**  
Children's Hospital of Philadelphia  
**Raymond Fenwick**  
Biotech Laboratories  
**Chris Fields**  
National Center for Genome Resources  
**Simon Foote**  
Walter and Eliza Hall Institute of  
Medical Research

**Richard Gibbs**  
Baylor College of Medicine  
**Eric Green**  
National Human Genome  
Research Institute, NIH  
**Richard Myers**  
Stanford University School of Medicine

**Phil Green**  
University of Washington  
**Kenshi Hayashi**  
Kyushu University  
**Philip Hieter**  
The Johns Hopkins University School  
of Medicine  
**Clare Huxley**  
St. Mary's Hospital Medical School  
**Howard J. Jacob**  
Medical College of Wisconsin  
**Alec Jeffreys**  
University of Leicester  
**Mark Johnston**  
Washington University School of  
Medicine  
**Mary-Claire King**  
University of Washington  
**Ben Koop**  
University of Victoria  
**Pui-Yan Kwok**  
Washington University School of  
Medicine  
**Ulf Landegren**  
Uppsala Biomedical Center  
**Mark Lathrop**  
The Wellcome Trust Centre  
**Michael Lovett**  
University of Texas Southwestern  
Medical Center  
**Jen-I Mau**  
Genome Therapeutics Corporation  
**Douglas Marchuk**  
Duke University Medical Center  
**Thomas Marr**  
Cold Spring Harbor Laboratory  
**W. Richard McCumbie**  
Cold Spring Harbor Laboratory  
**Susan Naylor**  
University of Texas Health Science  
Center  
**David Nelson**  
Baylor College of Medicine

## Reviews Editor

**Allison Stewart**  
Cambridge, UK

STEENBOCK  
MEMORIAL LIBRARY

APR 13 1998

**Maynard Olson**  
University of Washington  
**Svante Pääbo**  
University of Munich  
**Leena Peltonen**  
National Public Health Institute,  
Helsinki  
**David Porteous**  
MRC Human Genetics Unit  
Western General Hospital, Edinburgh  
**Roger Reeves**  
Johns Hopkins University School of  
Medicine  
**Bruce Roe**  
University of Oklahoma  
**Rodney Rothstein**  
Columbia University College of P&S  
**Gerald Rubin**  
University of California, Berkeley  
**Lloyd Smith**  
University of Wisconsin-Madison  
**Randall Smith**  
Baylor College of Medicine  
**Marcelo Bento Soares**  
University of Iowa  
**William Studier**  
Brookhaven National Laboratory  
**Grant Sutherland**  
Women's and Children's Hospital,  
Adelaide  
**Barbara Trask**  
University of Washington  
**Gert-Jan B. van Ommen**  
Leiden University  
**Robert D. Wells**  
University of Utah  
**Jean Weissenbach**  
Genethon, CNRS  
**Richard Wilson**  
Washington University School of  
Medicine  
**James Womack**  
Texas A&M University

## Editorial Office

Cold Spring Harbor Laboratory Press  
1 Bungtown Road  
Cold Spring Harbor, New York 11724  
Phone (516) 367-6834  
Fax (516) 367-8334  
<http://www.cshl.org>

## Editorial/Production

**Nadine Dumser**, Technical Editor  
**Kristin Kraus**, Production Editor  
**Cynthia Grimm**, Production Editor  
**Peggy Calicchia**, Editorial Secretary

## IN THE UNITED STATES PATENT AND TRADEMARK OFFICE



In re application of:

Li *et al.*

Appl. No. 09/720,086

102(e): July 23, 2001

For: *De Novo* DNA Cytosine  
Methyltransferase Genes,  
Polypeptides and Uses Thereof

Confirmation No.: 6968

Art Unit: 1642

Examiner: Harris, A. M.

Atty. Docket: 0609.4560002/KRM/DJN

## Declaration Under 37 C.F.R. § 1.132 of Kenneth D. Bloch, M.D.

Commissioner for Patents  
P.O. Box 1450  
Alexandria, VA 22313-1450

Sir:

I, the undersigned, Kenneth D. Bloch, M.D., residing at 80 Park Street, Apt. 32, Brookline, Massachusetts, 02446, declare and state as follows:

1. I am currently employed by Massachusetts General Hospital in the Cardiovascular Research Center where I conduct and supervise research in the field of molecular cardiology. I worked with Dr. En Li, co-inventor of the captioned application, in the Cardiovascular Research Center for 10 years until 2003 when Dr. Li moved to Novartis. I am also an Associate Professor of Anesthesiology and Medicine at Harvard Medical School and have extensive experience in molecular biology and DNA cloning and sequencing.
2. A current *curriculum vitae* is appended hereto as EXHIBIT A.
3. I have been informed that human DNMT3A cDNA clone, represented in the captioned patent application as SEQ ID NO:3, was deposited with the ATCC on July 10, 1998, and assigned ATCC Deposit No. 98809. I have also been informed that the deposit date of July 10, 1998 was prior to the filing date of the second provisional application, App. No. 60/093,993, filed July 24, 1998, the benefit of which is claimed. Finally, I have been informed that the '993 application includes the sequence information and references the deposit of the sequenced material on page 16, lines 1-2, of the specification.
4. In November 2004, Applicants had samples withdrawn of the human DNMT3A cDNA clone contained within ATCC Deposit No. 98809. At the Applicants request, I have sequenced the nucleotides that span the coding region of the deposited human DNMT3A cDNA clone contained in ATCC Deposit No. 98809. A nucleotide alignment that spans the coding regions of the sequenced human DNMT3A cDNA clone contained in

- 2 -

Li et al.  
Appl. No. 09/720,086

ATCC Deposit No. 98809 and currently amended **SEQ ID NO:3** is shown in Fig. 1 of EXHIBIT B.

5. The amendment to the sequence listing, which was filed on July 23, 2001, corrected six nucleotide errors in the coding sequence of **SEQ ID NO:3** (see bolded nucleotides at nucleotide positions 940, 1476, 1479, 1570, 2024 and 2119 of amended **SEQ ID NO:3** in Fig. 1 of EXHIBIT B). The amendment also deleted original nucleotides 1-123 of **SEQ ID NO:3**, which does not include any DNMT3A coding sequence.

6. The deposited clone recited in ¶¶4 and 5, above (i.e., ATCC Deposit No. 98809) is the same as the deposited clone recited in the above-captioned application. The six nucleotides in the coding region of **SEQ ID NO:3** that were corrected by the amendment of July 23, 2001 correspond to the sequence contained in ATCC Deposit No. 98809. It is well known that sequencing errors are a common problem in Molecular Biology. Peter Richterich, Estimation of Errors in 'Raw' DNA Sequences: A Validation Study, 8 *Genome Research* 251-59 (1998)(EXHIBIT C). I believe that one skilled in the art would have sequenced the deposited material and recognized the sequencing errors.

7. My sequencing of ATCC Deposit No. 98809 also revealed that nucleotides 539-584 within the coding region of amended **SEQ ID NO:3** are deleted in the deposited cDNA. The deletion causes a frame shift in the coding region of the deposited cDNA and predicts a truncated protein of 145 amino acids. An amino acid alignment of the predicted amino acid sequence encoded by the human DNMT3A cDNA clone contained in ATCC Deposit No. 98809 and the predicted amino acid sequence encoded by amended **SEQ ID NO:3** is shown in Fig. 2 of EXHIBIT B. The bolded sequence in Fig. 2 corresponds to the predicted encoded amino acid sequence downstream of the nucleotide deletion in ATCC Deposit No. 98809 and represents a point of divergence compared with the predicted amino acids encoded by currently amended **SEQ ID NO:3**.

8. Currently amended **SEQ ID NO:3** and **SEQ ID NO:3** as originally filed in U.S. Appl. Nos. 60/090,906 and 60/093,993, to which priority is claimed, do not harbor the deletion, and encode a protein having 912 amino acids that is homologous to mouse Dnmt3a.

9. Like DNA sequence errors, it is known that errors in DNA cloning may occur in molecular biology. Deletion errors may occur and may be caused by, *inter alia*, inadvertent digestion of DNA by restriction endonucleases or exonucleases, or by recombination events during propagation of the DNA in bacterial hosts. The deletion at nucleotides 539-584 in **SEQ ID NO:3** present in ATCC Deposit No. 98809 is an obvious error. I believe that one skilled in the art would have sequenced the deposited material and recognized the deletion as an error. My belief is based upon the following: First, the deletion found in ATCC Deposit No. 98809 is not present in **SEQ ID NO:3** as originally filed or as amended. Second, the deletion is not present in the DNA sequence of the closely related mouse homolog, **SEQ ID NO:1**. Third, the deletion causes a frame shift in the reading frame of **SEQ ID NO:3** and predicts a truncated protein product compared with that

- 3 -

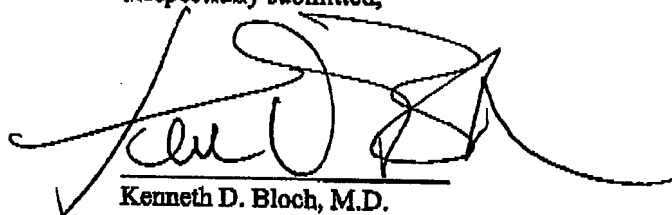
Li et al.  
Appl. No. 09/720,086

encoded by SEQ ID NO:3 as originally filed and as amended. Fourth, the amino acids encoded by the nucleotide sequence downstream of the deletion bear no similarity to the amino acids encoded by SEQ ID NO:3 or the mouse homolog of Dnmt3a, encoded by SEQ ID NO:1. Finally, an examination of the sequence reveals two large open reading frames (ORF) in the sequence in different frames. See Fig. 3 of EXHIBIT B. The ORFs correspond to the amino acid residues of DNMT3A upstream and downstream of the deletion. The presence of two large ORFs in different frames indicates a possible frameshifting sequence error or deletion. All of these factors indicate that the deletion present in ATCC Deposit No. 98809 is an error, and would be recognized as such by a person of ordinary skill in the art.

10. It is my belief that the combination of ATCC Deposit No. 98809, which discloses the six nucleotides in the coding region of SEQ ID NO:3 amended on July 23, 2001, in combination with SEQ ID NO:3 as originally filed in U.S. Appl. Nos. 60/090,906 and 60/093,993, which disclose nucleotides 539-584 of amended SEQ ID NO:3, clearly conveys to someone skilled in the art the entire nucleotide sequence of amended SEQ ID NO:3.

11. I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the present patent application or any patent issued thereon.

Respectfully submitted,



Kenneth D. Bloch, M.D.

Date: 11/7/2005

# EXHIBIT A



## **CURRICULUM VITAE**

### **PART I: General Information**

**DATE PREPARED:** July 5, 2005

**Name:** Kenneth D. Bloch

**Office Address:** Cardiovascular Research Center  
Massachusetts General Hospital  
149 13<sup>th</sup> Street (149-4201)  
Charlestown, MA 02129  
(617) 724-9540

**Home Address:** 80 Park Street, Apartment 32, Brookline, MA 02446

**E-Mail:** kdbloch@partners.org

**FAX:** (617) 726-5806

**Place of Birth:** New York, NY

#### **Education:**

1978	Sc.B.	Brown University (Biomedicine)
1981	M.D.	Brown University

#### **Postdoctoral Training:**

##### **Internship and Residencies:**

1981-1984	Internal Medicine, Massachusetts General Hospital
-----------	---

##### **Fellowships:**

1981-1984	Clinical Fellow in Medicine, Harvard Medical School
1984-1987	Research Fellow in Genetics, Harvard Medical School
1987-1989	Clinical and Research Fellow in Medicine, Harvard Medical School

#### **Licensures and Certification:**

1984-	Commonwealth of Massachusetts, Registered Physician,
1985	Diplomate, American Board of Internal Medicine
1989	Diplomate, Subspecialty Board of Cardiovascular Diseases, American Board of Internal Medicine

#### **Academic Appointments:**

1989-1990	Instructor in Medicine, Harvard Medical School
1990-1997	Assistant Professor of Medicine, Harvard Medical School

1997- Associate Professor of Medicine, Harvard Medical School  
 2005- Associate Professor of Anaesthesia, Harvard Medical School

**Hospital or Affiliated Institution Appointments:**

1990-1996 Assistant in Medicine, Massachusetts General Hospital  
 1996-1999 Assistant Physician, Massachusetts General Hospital  
 1999-2005 Associate Physician, Massachusetts General Hospital  
 2005- Physician, Massachusetts General Hospital

**Other Professional Positions and Major Visiting Appointments:** none

**Hospital and Health Care Organization Service Appointments:**

1990- Attending Physician, Cardiology and Medical Services,  
 Massachusetts General Hospital  
 1991- Practice Member, Cardiac Unit Associates,  
 Massachusetts General Hospital

**Major Administrative Responsibilities:**

1989-1992 Principal Investigator, Cardiac Unit, Massachusetts General Hospital  
 1990- Associate Program Director, Fellowship Program in Cardiovascular  
 Disease, Massachusetts General Hospital  
 1992- Principal Investigator, Cardiovascular Research Center, Massachusetts  
 General Hospital  
 1993-1997 Preceptor, Training Grant to the Cardiovascular Research Center (PI:  
 Mark C. Fishman), Massachusetts General Hospital  
 1997-2002 Co-Investigator, Training Grant to the Cardiovascular Research Center  
 (PI: Mark C. Fishman), Massachusetts General Hospital  
 2002-2004 Interim Director, Cardiovascular Research Center, Massachusetts  
 General Hospital  
 2002- Principal Investigator, Training Grant to the Cardiovascular Research  
 Center, Massachusetts General Hospital

**Major Committee Assignments:**

Hospital:

1991 Member, Selection Committee for the Fellowship Program in  
 Cardiovascular Disease, Massachusetts General Hospital  
 1992-1996 Member, Subcommittee on Review of Research Proposals, Committee  
 on Research, Massachusetts General Hospital  
 1997- Member, Steering Committee coordinating the integration of the  
 clinical cardiology fellowship programs at Brigham and Women's  
 Hospital and Massachusetts General Hospital

Regional:  
1991-1994

2002

Member, Research Peer Review Committee, American Heart Association, Massachusetts Affiliate, Inc.  
Member, Northeast Peer Review Study Group for Lipids, Thrombosis & Vascular Wall Biology, American Heart Association

National:  
1992-1995

1996-1997

2000

2000-2002

2002-

2002-

2005

2005-

2005-

Member, Molecular Signaling Research Study Committee, American Heart Association, National Center  
Member, Lung Biology and Pathology Study Section, National Heart Lung and Blood Institute  
Member-At-Large, Executive Committee of the Council on Cardiopulmonary and Critical Care, American Heart Association, National Center  
Member, Program Committee, Council on Cardiopulmonary and Critical Care, American Heart Association, National Center  
Chairman, Program Committee, Council on Cardiopulmonary and Critical Care, American Heart Association, National Center  
Member, Committee on Scientific Sessions Program, American Heart Association  
Member, Respiratory Integrative Biology and Translational Research Study Section, National Heart Lung and Blood Institute  
Member, Research Committee, American Heart Association  
Member, Future of Scientific Sessions Task Force, American Heart Association

**Professional Societies:**

1997-

2000

2001

American Heart Association, Council on Basic Cardiovascular Science  
American Heart Association, Council on Cardiopulmonary and Critical Care  
American Society for Clinical Investigation

**Editorial Boards:**

1989-

Ad hoc reviewer:  
American Journal of Pathology  
American Journal of Physiology  
American Journal of Respiratory Cell and Molecular Biology  
Anesthesiology  
Circulation  
Circulation Research  
Journal of Applied Physiology  
Journal of Biological Chemistry  
Journal of Clinical Investigation  
New England Journal of Medicine  
Nature Medicine

## Proceedings of the National Academy of Sciences

### Awards and Honors:

1978	Magna cum laude, Brown University
1981	Patricia McCormick Memorial Prize, Brown University
1984-1986	Research Fellowship, Leukemia Society of America
1986-1987	Postdoctoral Fellowship, Pfizer Pharmaceuticals
1996-1998	Grant-in-Aid, American Heart Association, National Center (declined)
1996	Finalist, Circulation Council Cardiovascular Research Competition
2000	Second Prize for the basic science research abstracts from the Massachusetts Thoracic Society
2001	Charter Fellow of the Council on Basic Cardiovascular Sciences, American Heart Association
2002	Inaugural Fellow of the Basic Sciences Council, American Heart Association
2003	Inaugural Fellow of the Council on Cardiopulmonary, Perioperative, and Critical Care, American Heart Association

### Part II: Research, Teaching, and Clinical Contributions

#### A. Narrative Report

Dr. Bloch's research has focused on cardiovascular biology and the molecular mechanisms regulating vascular tone and ventricular remodeling. Dr. Bloch established his own laboratory at MGH in 1991 and became a principal investigator in the Cardiovascular Research Center in 1992.

When he established his laboratory, it appeared that there was a single constitutively-expressed nitric oxide (NO) synthase (NOS) in brain and endothelium. Dr. Bloch's group cloned the NOS isoform responsible for endothelial NO production, NOS3. Dr. Bloch's group discovered that NOS3-deficient mice are more susceptible to hypoxia-induced pulmonary vascular remodeling. Enhanced pulmonary NO (either by inhalation of NO gas or by aerosol delivery of an adenovirus specifying NOS attenuates pulmonary vasoconstriction and prevents this pulmonary vascular remodeling. Moreover, they found that NO inhalation could prevent pulmonary vascular remodeling even before the development of pulmonary vasoconstriction. Taken together, these studies have important implications for the treatment of children with congenital heart disease, in whom pulmonary vascular remodeling precedes the development of pulmonary hypertension (which is often fatal). Dr. Bloch and his colleagues were the first to show that NO inhalation also has systemic vascular effects including attenuation vascular neointima formation after balloon injury and prevention re-thrombosis after coronary artery thrombolysis (the latter is markedly potentiated by inhibitors of cGMP-metabolizing phosphodiesterases). These innovative studies have direct clinical implications for the care of patients with acute coronary syndromes.

Dr. Bloch has brought into focus the fact that regulation of NO responsiveness may be as important as regulation of NO production. Dr. Bloch found that prolonged exposure of vascular smooth muscle cells (VSMC) to NO or pro-inflammatory cytokines decreases function of soluble guanylate cyclase (sGC; the enzyme responsible for cGMP synthesis in response to NO), desensitizing the cells to NO. Dr. Bloch also showed that cGMP-dependent protein kinase, an enzyme responsible for vasodilation, also has a critical role in determining the sensitivity of these cells to the antiproliferative and proapoptotic effects of NO and cGMP.

Dr. Bloch's group has contributed importantly to a second research area--the structure and function of the "nuclear body", a nuclear structure that appears to have a critical role in oncogenesis, gene transcription, and the cellular response to viral infection. They have cloned two new nuclear body components, both of which appear to be transcription factors and one of which is a co-activator of nuclear hormone receptors.

Dr. Bloch's research group has elucidated an important role for NO synthesized by endothelial nitric oxide synthase (NOS3) in the left ventricular remodeling induced by a variety of hemodynamic challenges. In addition, they have explored the mechanisms responsible for the impairment of hypoxic pulmonary vasoconstriction associated with pulmonary injury associated with endotoxemia and volutrauma.

Most recently, Dr. Bloch's group has begun a new line of research directed at understanding how mutations in the gene encoding bone morphogenetic protein receptor type 2 (BMPR2) cause primary pulmonary hypertension. They have observed that mice carrying one copy of a mutant BMPR2 gene have mild pulmonary hypertension associated with abnormalities of pulmonary vascular structure.

Dr. Bloch has also made important contributions in the clinical and research training of cardiology fellows. He is a primary research mentor for the MGH cardiology fellows guiding them to, and supporting them in, the best research training opportunities HMS has to offer. In the past five years, Dr. Bloch has played a pivotal role in the integration of the cardiology fellowship programs at the Brigham and Women's Hospital and MGH providing fellows with exposure to outstanding clinical experiences at both institutions. Since 2002, Dr. Bloch has served as principal investigator of the T32 training grant awarded to the Cardiovascular Research Center. In this role, Dr. Bloch supervises the mentoring and career development of 10 post-doctoral cardiovascular scientists each year. From 2002 through 2004, Dr. Bloch served as the Interim Director of the CVRC fostering the creativity and productivity of ten faculty members and more than 40 post-doctoral research fellows.

#### **B. Funding Information (Research):**

##### Past:

- |           |  |
|-----------|--|
| 1991-1992 | Research Grant, Dr. Louis Skarow Memorial Fund (PI: K. Bloch)<br>"Gene expression in a model of pulmonary hypertension." |
| 1991-1994 | Grant-in-Aid, American Heart Association, National Center<br>(PI: K. Bloch) "Pulmonary expression of endothelin genes."  |

1991-1996 NHLBI/R29 (PI: K. Bloch)  
 "Biosynthesis of the endothelin family of vasoactive peptides."  
 1995-1996 Sponsored Research Agreement through the Cardiovascular  
 Research Center from Bristol-Myers Squibb (PI: K. Bloch)  
 "Isolation and characterization of novel vascular nitric oxide  
 synthases."  
 1996-1997 NHLBI/R01 (Co-PI: K. Bloch)  
 "The pulmonary response to inhaled particulates."  
 1996-2000 NHLBI/R01 (PI: K. Bloch)  
 "Nitric oxide/cGMP signal transduction in pulmonary injury."  
 1996-2001 Established Investigator, American Heart Association, National Center  
 (PI: K. Bloch) "Regulation of a nitric oxide receptor component, the  $\beta 1$   
 subunit of soluble guanylate cyclase."  
 1998-2002 NHLBI/T32 (Co-PI: K. Bloch; PI: M.C. Fishman)  
 "Cell and molecular training for cardiovascular biology."  
 1998-2003 NHLBI/R01 (PI: K. Bloch)  
 "Nitric oxide/cGMP signal transduction in vascular injury."  
 2001-2002 Pfizer Pharmaceuticals (Co-PI: K. Bloch; PI: M.J. Semigran)  
 "Evaluation of the effects of sildenafil with and without inhaled nitric  
 oxide (NO) on platelet-mediated thrombosis and cardiac function in a  
 canine model of cyclic coronary artery occlusion."

Current:

1996- NHLBI/R01 (Co-PI: K. Bloch; PI: W. Zapol)  
 "Studies of inhaled nitric oxide."  
 2002- INO Therapeutics, Inc. (PI: K Bloch)  
 "Laboratory-based initiatives for the further development of the  
 therapeutic potential of inhaled nitric oxide."  
 2002- NHLBI/T32 (PI: K. Bloch)  
 "Cell and molecular training for cardiovascular biology."  
 2003- NHLBI/R01 (PI: K. Bloch)  
 "Nitric oxide synthase 3 and left ventricular remodeling"  
 2003- NHLBI/R01 (PI: K. Bloch)  
 "BMPR2 in the pathogenesis of pulmonary hypertension"

**C. Report of Current Research Activities (Bench and Clinical Research):**

Project 1: Studies of inhaled nitric oxide. (Co-PI: K. Bloch)  
 Project 2: Evaluation of the systemic effects of breathing nitric oxide (PI: K. Bloch)  
 Project 3: Role of nitric oxide in left ventricular remodeling (PI: K. Bloch)  
 Project 4: Role of BMPR2 in the pathogenesis of pulmonary hypertension (PI: K. Bloch)  
 Project 5: Evaluation of nitric oxide inhalation to treat cardiogenic shock complicating  
 right ventricular infarction. (Co-PI: K. Bloch)

## **D. Report of Teaching**

### **1. Local Contributions**

#### **a. medical school**

- 1977 Brown University, Biomed 110, Biophysics  
course director: Babette Stewart  
Teaching Assistant  
50 undergraduates (approx.)  
3-5 hours preparation and contact time/week (approx.)  
semester course
- 1978-1981 Brown University, Biomed 6, Introduction to Physiology  
course director: Peter Stewart  
Teaching Assistant  
50 undergraduates (approx.)  
3-5 hours preparation and contact time/week (approx.)  
semester course
- 1986 Harvard University, Genetics 700.0, Fundamentals of Genetics  
course director: Philip Leder  
Teaching Assistant  
15 medical students (approx.)  
3 hours preparation and contact time/week (approx.)  
semester course
- 1990-1992 Harvard Medical School, Introduction to Clinical Medicine  
Clinical Mentor  
2-3 medical students, 3 sessions/week, 2 hours/session, 3 weeks/year  
total: 9 hours preparation time, 18 hours contact time
- 1993-1996 Harvard Medical School, Patient-Doctor II Course  
course directors: Katherine Treadway and Diane Fingold  
Preceptor  
50 medical students (approx.)  
3 sessions, 2 hours/session, 1/year  
total: 9 hours preparation and contact time

#### **b. graduate medical course:** none

#### **c. local invited teaching presentations**

- 1993-1994 "Ethical Conduct of Research" Course, Massachusetts General  
Hospital, Preceptor, 5-10 postdoctoral fellows, 2 sessions/year,  
2 hours/session, prep time: 1 hour/session.
- 1994 Cardiology Grand Rounds, Massachusetts General Hospital, Lecturer;

- 1996 50 attendees: medical students, residents, clinical and research fellows, faculty; 4 hrs prep and 1 hr contact time.  
Anesthesiology Grand Rounds, Massachusetts General Hospital; Lecturer; 50 attendees: medical students, residents, clinical and research fellows, faculty; 4 hrs prep and 1 hr contact time.
- 1995 Wellman Laboratories of Photomedicine Symposium, Massachusetts General Hospital, Session Chair; 50 attendees: medical students, residents, clinical and research fellows, faculty; 1 hr prep and contact time.
- 1997 Cardiology Grand Rounds, Massachusetts General Hospital, Lecturer; 50 attendees: medical students, residents, clinical and research fellows, faculty; 4 hrs prep and 1 hr contact time.
- 1998 Clinical Fellows' Core Curriculum Lecture Series, Massachusetts General Hospital, Lecturer; 5 attendees: medical students and fellows; 4 hrs prep time and 1 hr contact time.
- 1999 Pulmonary and Critical Care Unit Research Conference, Massachusetts General Hospital, Lecturer; 50 attendees: medical students, residents, clinical and research fellows, faculty; 4 hrs prep and 1 hr contact time.
- 2000 Center for the Prevention of Cardiovascular Disease, Department of Nutrition, Harvard School of Public Health, Lecturer; 50 attendees: medical students, residents, clinical and research fellows, faculty; 4 hrs prep and 1 hr contact time.
- 2001 West Roxbury Veterans Administration Hospital, Cardiology Grand Rounds; 30 attendees: medical students, residents, clinical and research fellows, faculty; 4 hrs prep and 2 hr contact time.
- 2003 Brigham and Women's Hospital, Vascular Research Seminar; 30 attendees: medical students, residents, clinical and research fellows, faculty; 4 hrs prep and 2 hr contact time.
- 2004 Massachusetts General Hospital, Cardiology Grand Rounds; 50 attendees: nurses, medical students, residents, clinical and research fellows, faculty; 5 hrs prep and 1 hr contact time.
- 2005 Massachusetts General Hospital, Critical Care Research Retreat; 50 attendees: clinical and research fellows, faculty; 10 hours prep time, 4 hours contact time, 10 minute lecture

d. continuing medical education courses: none

e. advisory and supervisory responsibilities

- 1989- Attending Physician, Private and Ward Medical Services, Massachusetts General Hospital (variable hours/year)
- 1989-1992 Principal Investigator, Cardiovascular Research Center, Massachusetts General Hospital, scientific supervisor for Research Fellows, including one Associate Professor and one Assistant Professor of Anesthesia (1,000 hours/year).
- 1990- Attending Physician, Cardiology Consult Service, Massachusetts General Hospital (140 hours/year).
- 1991- Research Advisor, Cardiology Fellowship Program, Massachusetts



- General Hospital, 25-35 Fellows per year in clinical and research training (200 hours/year).
- 1992-2002      Attending Physician, Coronary Care Unit, Massachusetts General Hospital (variable hours/year).
- 1992-      Principal Investigator, Cardiovascular Research Center, Massachusetts General Hospital, scientific supervisor for Research Fellows, including two Assistant Professors of Anesthesia (1,000 hours/year).
- 2002-      Principal Investigator, NIH-sponsored program (T32) to train 10 post-doctoral cardiovascular scientists each year (200 hours/year)

f. teaching leadership role

- 1990-1999      Cardiac Unit Research Seminar Series, Massachusetts General Hospital; Organizer; A weekly series of presentation by staff and senior fellows designed to highlight research in the Cardiac Unit.
- 1995-1996      Cardiac Unit Society of Fellows, Massachusetts General Hospital; Organizer; A quarterly series of symposia presented by research fellows at the home of the Chief of the MGH Cardiac Unit.
- 1997      Society of Cardiology Fellows, Massachusetts General Hospital and Brigham and Women's Hospital; Co-Organizer; A quarterly series of symposia designed to foster scientific communication and collaboration between MGH and BWH and to highlight research opportunities for fellows in cardiology training.
- 2002-      CVRC Seminar Series; Organizer; A weekly seminar series presented by visiting scientists in the MGH Cardiovascular Research Center

g. names of advisees and trainees/current positions

- 1989-1990      Charles C. Hong, MD/PhD, Clinical and Research Fellow, Cardiology Division, Massachusetts General Hospital
- 1990-1991      Robert Schott, MD, MPH, private practice
- 1991-1993      Akito Shimouchi, MD, Assistant Professor of Medicine, National Cardiovascular Center Research Institute, Osaka, Japan
- 1990-1991      Stefan P. Janssens, MD, PhD, Associate Professor of Medicine, Cardiac Unit and Laboratory for Molecular and Vascular Biology, University Hospital Gasthuisberg, Leuven, Belgium
- 1992-1994, 2001-      Noriko Kawai, MD, PhD, Research Fellow in Medicine, Cardiovascular Research Center, Massachusetts General Hospital
- 1992-1993      John J. Lepore, MD, Instructor of Medicine, University of Pennsylvania Medical Center, Department of Medicine, Cardiovascular Medicine Division and Molecular Cardiology Laboratory
- 1993-1994      Johanna Wolfram, MD, University Clinic for Internal Medicine II, Department of Cardiology, General Hospital, Vienna, Austria
- 1993-1999      Lucienne Sanchez, MD, Instructor in Pediatrics, Harvard Medical School/Massachusetts General Hospital

1993-2004	Jesse D. Roberts, MD, Associate Professor of Anesthesiology in Pediatrics, Harvard Medical School/Massachusetts General Hospital
1994-1995	Jeffrey Thomas, MD, Associate Professor of Neurosurgery, Chief of Neurovascular and Neurointerventional Surgery, Division of Neurosurgery, The University of New Mexico Health Sciences Center
1994-1996	Alexandra Holzmann, MD, Department of Anesthesiology, University of Heidelberg
1995-2004	Heling Liu, MD, on leave to care for her children
1995-1998	Jean-Daniel Chiche, MD, Professor, Cochin University, Paris, France
1996-1997	Anita Honkanen, MD, private practice
1996-1998	Masao Takata, MD/PhD, Assistant Professor, Department of Anaesthetics and Intensive Care, Imperial College School of Medicine, Hammersmith Hospital, London, United Kingdom
1996-1998	Douglas Wirthlin, MD, Assistant Professor of Surgery, Vascular Surgery, University of Alabama at Birmingham
1996-1998	Joerg Weimann, MD, Professor of Anesthesia, Department of Anaesthesiology and Intensive Care Medicine, Charité -Berlin Medical School, Campus Benjamin Franklin
1998-2002	Galina Filippov, MD, Research Scientist, Omnigene Bioproducts Inc.
1998, 2000-2001	Zena Quezado, MD, Chief, Department of Anesthesia and Surgical Services, Warren G. Magnuson Clinical Center, National Institutes of Health
1998-2000	Roman Ullrich, MD, Associate Professor of Anesthesia and Intensive Care Medicine, Vienna General Hospital, Medical University of Vienna
1999-2001	Hiroshi Nakajima, MD, Neurosurgery Residency, Tokyo Women's Medical College, Tokyo, Japan
1999-2001	Pini Orbach, PhD, Project Manager, Drug Development, Perdix Pharmaceuticals, Inc.
1999-	Fumito Ichinose, MD, PhD, Assistant Professor of Anesthesia, Harvard Medical School/Massachusetts General Hospital
2001-2003	Aimee Limbach, PhD, Post-doctoral Fellow, Center for Human Molecular Genetics, Munroe-Meyer Institute, University of Nebraska Medical Center
2002-2003	Elisabeth Choe, MD, Resident in Internal Medicine, University of Texas Southwestern
2002-2004	Cornelius Busch, MD, Resident in Anesthesia, Department of Anesthesiology, University of Heidelberg
2003	Claire Mayeur, MD, Resident in Anesthesiology, Lille, France
2003-	Hideyuki Beppu, MD, PhD, Instructor in Medicine, Cardiovascular Research Center, Massachusetts General Hospital
2003-	Manu Buys, PhD, Research Fellow in Medicine, Cardiovascular Research Center, Massachusetts General Hospital

- 2003- Paul Yu, MD, PhD, Clinical and Research Fellow in Medicine,  
Cardiovascular Research Center, Massachusetts General Hospital
- 2004-2005 David Bayne, undergraduate student, Harvard University
- 2004- Ryuji Hataishi, MD, PhD, Research Fellow in Anesthesia,  
Massachusetts General Hospital
- 2004- Rajeev Malhotra, MD, Resident in Internal Medicine, Massachusetts  
General Hospital
- 2004- Tomas Neilan, MD, Clinical and Research Fellow in Medicine,  
Cardiac Ultrasound Laboratory and Cardiovascular Research Center,  
Massachusetts General Hospital
- 2005- Sarah Blake, MD, Research Fellow in Anesthesia, Massachusetts  
General Hospital

## 2. Regional, National, or International Contributions

- 1994 Invited Lecturer, American College of Cardiology, Dallas, TX
- 1994 Invited Lecturer, Pfizer Pharmaceuticals, Groton, CT
- 1994 Invited Lecturer, University of Leuven, Belgium
- 1995 Invited Lecturer, St. Elizabeth's Hospital, Cardiovascular Research  
Seminar, Boston, MA
- 1996 Invited Lecturer, Boston Heart Foundation, Boston, MA
- 1996 Invited Lecturer, Georgia Medical College, Vascular Biology  
Division, Atlanta, GA
- 1997 Invited Lecturer, Harvard Medical School, Vascular Biology  
Seminar, Boston, MA
- 1997 Invited Lecturer, Boston University, Whittaker Foundation, Boston, MA
- 1998 Invited Lecturer, Oregon Health Sciences University, Cardiology  
Division, Oregon
- 1998 Invited Lecturer, Brigham and Women's Hospital, Cardiology  
Division, Monday Morning Research Conference, Boston, MA
- 1998 Invited Lecturer, New York Medical College, Department of  
Pharmacology Seminar, New York, NY
- 1999 Invited Lecturer, Tufts University School of Medicine, Dept. of  
Medicine, Boston, MA
- 1999 Invited Lecturer, Tufts University School of Medicine/New England  
Medical Center, Pulmonary and Critical Care Division, Boston, MA
- 1999 Invited Lecturer, Millenium Pharmaceuticals, Inc., Cambridge, MA
- 1999 Invited Lecturer, University of Washington, Cardiology Dept.,  
Seattle, WA
- 1999 Invited Lecturer, University of Alabama at Birmingham, Dept. of  
Pathology, Birmingham, AL
- 1999 Invited Lecturer, 3<sup>rd</sup> International Society for Medical Gases  
Meeting, Heidelberg, Germany
- 1999 Invited Lecturer, National Institute of Health, Critical Care  
Medicine, Bethesda, MD
- 2004 Invited Lecturer, INO Therapeutics Inc. Scientific Advisory Board,  
Chatham, MA

- 2002      Invited Lecturer, Medical College of Wisconsin, Milwaukee, Wisconsin
- 2002      Invited Lecturer, American Heart Association Scientific Sessions, Chicago, IL
- 2003      Invited Lecturer, Department for Molecular and Biomedical Research, Universiteit Gent, Belgium
- 2003      Invited Lecturer, Whitaker Cardiovascular Institute, Boston University Medical School
- 2003      Invited Lecturer, Cardiology Division, University of Alberta, Edmonton, Canada
- 2004      Invited Lecturer, American Heart Association, Northeast Affiliate, Symposium: Launching a Career in Cardiovascular Research
- 2004      Invited Lecturer, Cardiology Grand Rounds, Dartmouth-Hitchcock Medical Center—"Nitric oxide synthases in ventricular remodeling: insights gained from genetically-modified mice."
- 2004      Invited Lecturer, Vascular Biology Seminar, Dartmouth-Hitchcock Medical Center—"Mechanisms regulating pulmonary vascular structure and function—roles of leukotrienes and bone morphogenetic proteins."
- 2004      Invited Lecturer, Department of Physiology Seminar, Louisiana State University-Shreveport—"Nitric oxide synthases in ventricular remodeling: insights gained from genetically-modified mice."
- 2004      Invited Lecturer, Cardiovascular Cell and Gene Therapy Conference II, Cambridge, MA—"Nitric oxide/cGMP signal transduction—implications for cardiovascular gene transfer."
- 2004      Invited Lecturer, Critical Therapeutics, Inc., Lexington, MA—"Mechanisms of pulmonary vascular dysfunction in lung injury: insights gained from genetically-modified mice."

#### **E. Report of Clinical Activities**

Dr. Bloch is a practicing cardiologist who maintains a practice within the Cardiac Unit Associates and Cardiology Division at the Massachusetts General Hospital. His practice consists of patients with cardiac problems of a moderate to high level of complexity referred to a tertiary care center.

### **Part III: Bibliography**

#### **Original Articles:**

1. Seidman CE, Bloch KD, Klein KA, Smith JA, Seidman JG. Nucleotide sequences of the human and mouse atrial natriuretic factor genes. *Science* 1984, 226:1206-1209.
2. Bloch KD, Scott JA, Zisfein JB, Fallon JT, Seidman CE, Matsueda GR, Margolies MN, Homcy CJ, Graham RM, Seidman JG. Biosynthesis and secretion of proatrial natriuretic factor by cultured rat cardiocytes. *Science* 1985, 230:1168-1172.
3. Graham RM, Bloch KD, Delaney VB, Bourke E, Seidman JG. Bartter's syndrome and the atrial natriuretic factor gene. *Hypertension* 1986, 8:549-551.
4. Ballerman BJ, Bloch KD, Seidman JG, Brenner BM. Atrial natriuretic peptide transcription, secretion, and glomerular receptor activity during mineralocorticoid escape. *J Clin Invest* 1986, 78:840-843.
5. Zisfein JB, Matsueda GR, Fallon JT, Bloch KD, Seidman CE, Seidman JG, Homcy CJ, Graham RM. Atrial natriuretic factor: assessment of its structure in atria and regulation of its biosynthesis with volume depletion. *J Mol Cell Cardiol* 1986, 18:917-929.
6. Bloch KD, Seidman JG, Naftilan JD, Fallon JT, Seidman CE. Neonatal atria and ventricles secrete atrial natriuretic factor via tissue-specific secretory pathways. *Cell* 1986, 47:695-702.
7. Bloch KD, Zisfein JB, Margolies MN, Homcy CJ, Seidman JG, Graham RM. Atrial natriuretic factor biosynthesis: a serum protease cleaves proANF into a 14-kilodalton peptide and ANF. *Am J Physiol* 1987, 252:E147-E151.
8. Bloch KD, Jones SW, Preibisch G, Seipke G, Seidman CE, Seidman JG. Proatrial natriuretic factor is phosphorylated by rat cardiocytes in culture. *J Biol Chem* 1987, 262:9956-9961.
9. Zeller R, Bloch KD, Williams BS, Arceci RJ, Seidman CE. Localized expression of the atrial natriuretic factor gene during cardiac embryogenesis. *Genes and Development* 1987, 1:693-698.
10. Mendez RE, Pfeffer JM, Ortola FV, Bloch KD, Anderson S, Seidman JG, Brenner BM. Atrial natriuretic peptide transcription, storage, and release in rats with myocardial infarction. *Am J Physiol* 1987, 253:H1449-H1455.
11. Lee RT, Bloch KD, Pfeffer JM, Pfeffer MA, Neer EJ, Seidman CE. Atrial natriuretic factor gene expression in ventricles of rats with spontaneous biventricular hypertrophy. *J Clin Invest* 1988, 81:431-434.

12. Bloch DB, Bloch KD, Iannuzzi M, Collins FS, Neer EJ, Seidman JG, Morton CC. The gene for the  $\alpha(i)1$  subunit of human G protein maps near the cystic fibrosis locus. *Am J Hum Gen* 1988, 42:884-888.
13. Seidman CE, Wong D, Bloch KD, Seidman JG. Cis-acting sequences that modulate atrial natriuretic factor gene expression. *Proc Natl Acad Sci USA* 1988, 85:4104-4108.
14. Kim S, Ang S-L, Bloch DB, Bloch KD, Kawahara Y, Tolman C, Lee R, Seidman JG, Neer EJ. Identification of cDNA encoding a new  $\alpha$  subunit of a human GTP-binding protein: expression of three  $\alpha(i)$  subtypes in human tissues and cell lines. *Proc Natl Acad Sci USA* 1988, 85:4153-4157.
15. Bloch KD, Zamir N, Seidman CE, Seidman JG. Ouabain induces secretion of proatrial natriuretic factor by neonatal rat atrial cardiocytes. *Am J Physiol* 1988, 255:E383-E387.
16. Ladenson PW, Bloch KD, Seidman JG. Modulation of atrial natriuretic factor (ANF) by thyroid hormone: mRNA and peptide levels in hypothyroid, euthyroid, and hyperthyroid rat atria and ventricles. *Endocrinology* 1988, 123:652-657.
17. Lee RT, Brock TA, Tolman C, Bloch KD, Seidman JG, Neer EJ. Subtype-specific increase in G-protein  $\alpha$ -subunit mRNA by interleukin  $1\beta$ . *FEBS Letters* 1989, 2:139-142.
18. Bloch KD, Friederich SP, Lee M-L, Eddy RL, Shows TB, Quertermous T. Structural organization and chromosomal assignment of the gene encoding endothelin. *J Biol Chem* 1989, 264:10851-10857.
19. Bloch KD, Eddy RL, Shows TB, Quertermous T. cDNA cloning and chromosomal assignment of the gene encoding endothelin 3. *J Biol Chem* 1989, 264:18156-18161.
20. Lee M-L, Bloch KD, Clifford JA, Quertermous T. Functional analysis of the endothelin-1 gene promoter: evidence for an endothelial cell-specific cis-acting sequence. *J Biol Chem* 1990, 265:10446-10450.
21. Lee M-L, de la Monte S, Ng S-C, Bloch KD, Quertermous T. Expression of the potent vasoconstrictor endothelin in the central nervous system. *J Clin Invest* 1990, 86:141-147.
22. Bloch KD, Hong CC, Eddy RL, Shows TB, and Quertermous T. cDNA cloning and chromosomal assignment of the endothelin 2 gene: vasoactive intestinal contractor peptide is rat endothelin 2. *Genomics* 1991, 10:236-242.
23. Janssens SP, Shimouchi A, Quertermous T, Bloch DB, and Bloch KD. Cloning and expression of a cDNA encoding human endothelium-derived relaxing factor/nitric oxide synthase. *J Biol Chem* 1992, 267:14519-14522.

24. Roberts JD, Chen T-Y, Kawai N, Wain J, Dupuy P, Shimouchi A, Bloch KD, Polaner D, Zapol WM. Inhaled nitric oxide reverses pulmonary vasoconstriction in the hypoxic and acidotic newborn lamb. *Circ Res* 1993, 72:246-254.
25. Shimouchi A, Janssens SP, Bloch DB, Zapol WM, Bloch KD. Cyclic AMP regulates soluble guanylate cyclase  $\beta$ 1 subunit gene expression in RFL-6 rat fetal lung fibroblasts. *Am J Physiol* 1993, 265:L456-L461.
26. Rovira I, Chen T-Y, Winkler M, Kawai N, Bloch KD, Zapol WM. Effects of inhaled nitric oxide on pulmonary hemodynamics and gas exchange in an ovine model of ARDS. *J Appl Physiol* 1994, 76:345-355.
27. Suen HC, Bloch KD, Donahoe PK. Antenatal glucocorticoid treatment corrects the pulmonary immaturity of congenital diaphragmatic hernia. *Pediatr Res* 1994, 35:523-529.
28. Cicila GT, Rapp JP, Bloch KD, Kurtz TW, Pravenec M, Kren V, Hong CC, Quertermous T, Ng S-C. Cosegregation of the endothelin-3, but not the endothelin-1, locus with blood pressure and relative heart weight in inbred Dahl rats. *J Hypertension* 1994, 12:643-651.
29. Horwitz MJ, Bloch KD, Kim NB, Amico JA. Expression of the endothelin 1 and oxytocin genes in the hypothalamus of the pregnant rat. *Brain Research* 1994, 648:59-64.
30. Bloch DB, Rabkina D, Quertermous T, Bloch KD. The immunoreactive region in a novel autoantigen contains a nuclear localization sequence. *Clin Immunol Immunopath* 1994, 72:380-389.
31. Roberts JD, Roberts CT, Jones RC, Zapol WM, Bloch KD. Nitric oxide inhalation reduces hypoxic pulmonary arterial remodeling, right ventricular hypertrophy, and growth retardation in the newborn rat. *Circ Res* 1995, 76:215-222.
32. de la Monte SM, Quertermous T, Hong CC, Bloch KD. Regional and maturation-associated expression of endothelin 2 in the gastrointestinal tract. *J Histochem Cytochem* 1995, 43:203-209.
33. Staples JF, Zapol WM, Bloch KD, Kawai N, Val VMF, Hochachka PW. Nitric oxide responses of air-breathing and water-breathing fish. *Am J Physiol* 1995, 268:R816-R819.
34. Kawai N, Bloch DB, Filippov G, Rapkina D, Suen H-C, Losty PD, Janssens SP, Zapol WM, de la Monte SM, Bloch KD. Constitutive endothelial nitric oxide synthase gene expression is regulated during lung development. *Am J Physiol* 1995, 268:L589-L595.

35. Bloch KD, Wolfram JR, Roberts JD, Zapol DG, Lepore JJ, Filippov G, Thomas JE, Brown D, Jacob HJ, Bloch DB. Three members of the nitric oxide synthase II gene family co-localize to human chromosome 17. *Genomics* 1995, 27:526-530.
36. Kurrek MM, Castillo L, Bloch KD, Tannenbaum SR, Zapol WM. Inhaled nitric oxide does not alter endotoxin-induced nitric oxide synthase activity during perfusion of the isolated rat lung. *J Appl Physiol* 1995, 79:1088-1092.
37. Bloch DB, Rabkina D, Bloch KD. The cell proliferation-associated antigen Ki-67 is a target of antibodies in the serum of MRL mice. *Lab Invest* 1995, 73:366-371.
38. Kurrek MM, Holzmann A, Filippov G, Winkler M, Zapol WM, Bloch KD. In vivo lipopolysaccharide pretreatment inhibits cGMP release from the isolated-perfused rat lung. *Am J Physiol* 1995, 269:L618-L624.
39. Huang PL, Huang ZH, Bloch KD, Moskowitz MA, Bevan JS, Fishman MC. Targeted disruption of the endothelial nitric oxide synthase gene causes hypertension. *Nature* 1995, 377:239-242.
40. Lee JS, Adrie C, Jacob HJ, Roberts JD, Zapol WM, Bloch KD. Chronic inhalation of nitric oxide inhibits neointima formation after balloon arterial injury in rats. *Circ Res* 1996, 78:337-342.
41. Janssens SP, Bloch KD, Nong Z, Gerard RD, Zoldhelyi P, Collen D. Adenovirus-mediated transfer of the human constitutive endothelial nitric oxide synthase gene to hypoxic rat lungs. *J Clin Invest* 1996, 98:317-324.
42. de la Monte SM, Bloch KD. Aberrant expression of the constitutive endothelial nitric oxide synthase gene in Alzheimer's disease. *Mol Chem Neuropath* 1997, 30:139-159.
43. Adrie C, Bloch KD, Moreno PR, Hurford WE, Guerrero L, Holt R, Zapol WM, Gold HK, Semigran MJ. Inhaled nitric oxide increases coronary artery patency after thrombolysis. *Circulation* 1996, 94:1919-1926.
44. Holzmann A, Bloch KD, Sanchez LS, Filippov G, Zapol WM. Hyporesponsiveness to inhaled NO in isolated-perfused lungs from endotoxin-challenged rats. *Am J Physiol* 1996, 271:L981-L986.
45. Bloch DB, de la Monte SM, Guigaouri P, Filippov A, Bloch KD. Identification and characterization of a leukocyte-specific component of the nuclear body. *J Biol Chem* 1996, 271:29198-29204.
46. Bloch KD, Filippov G, Sanchez LS, Nakane M, de la Monte SM. Pulmonary soluble guanylate cyclase, a nitric oxide receptor, is increased during the perinatal period. *Am J Physiol* 1997, 272:L400-L406.



47. Liu HL, Force T, Bloch KD. Nerve growth factor decreases soluble guanylate cyclase in rat pheochromocytoma PC12 cells. *J Biol Chem* 1997, 272:6038-6043.
48. Liu H-W, Anand A, Bloch KD, Christiani D, Kradin R. Expression of inducible nitric oxide synthase by macrophages in rat lung. *Amer Rev Resp Dis* 1997, 156:223-8.
49. Filippov G, Bloch DB, Bloch KD. Nitric oxide decreases stability of mRNAs encoding soluble guanylate cyclase subunits in rat pulmonary artery smooth muscle cells. *J Clin Invest* 1997, 100:942-948.
50. Head CA, Brugnara C, Martinez-Ruiz R, Kacmarek RM, Bridges KR, Kuter D, Bloch KD, Zapol WM. Low concentrations of nitric oxide increase oxygen affinity of sickle erythrocytes in vitro and in vivo. *J. Clin. Invest.* 1997, 100:1193-1198.
51. Sanchez LS, Filippov G, Zapol WM, Jones RC, Bloch KD. cGMP-binding, cGMP-specific phosphodiesterase gene expression is regulated during lung development. *Pediatr. Res.* 1998, 43:163-168.
52. Ichinose F, Adrie C, Hurford WE, Bloch KD, Zapol WM. Aerosolized zaprinast potentiates and prolongs the pulmonary vasodilation induced by breathing nitric oxide. *Anaesthesiol.* 1998, 88:410-416.
53. Steudel W, Scherrer-Crosbie M, Bloch KD, Weimann J, Huang PL, Picard MH, Zapol WM. Sustained pulmonary hypertension and right ventricular hypertrophy after chronic hypoxia in mice with congenital deficiency of nitric oxide synthase 3. *J. Clin. Invest.* 1998, 101:2468-2477.
54. Weimann J, Bauer H, Bigatello L, Bloch KD, Martin E, Zapol WM. ABO blood group and inhaled nitric oxide in acute respiratory distress syndrome. *Lancet* 1998, 351:1786-1787.
55. Chiche J-D, Schlutsmeyer SM, Bloch DB, de la Monte SM, Roberts JD, Filippov G, Janssens SP, Rosenzweig A, Bloch KD. Adenovirus-mediated gene transfer of cGMP-dependent protein kinase increases the sensitivity of cultured vascular smooth muscle cells to the antiproliferative and pro-apoptotic effects of nitric oxide/cGMP. *J. Biol. Chem.* 1998, 273:34263-34271.
56. Koike G, Chiche J-D, Shiozawa M, Simon JS, Szpirer J, Jacob HJ, Szpirer C, Bloch KD. Localization of rat genes in the nitric oxide signaling pathway: candidates for the pathogenesis of complex diseases. *Mamm. Genome.* 1999, 10:71-3.
57. Tyler RC, Muramatsu M, Abman SH, Stelzner TJ, Rodman DM, Bloch KD, McMurtry IF. Variable expression of endothelial NO synthase in three forms of rat pulmonary hypertension. *Am. J. Physiol.* 1999, 276:L297-303

58. Sohn YK, Ganju N, Bloch KD, Wands JR, de la Monte SM. Neuritic sprouting with aberrant expression of the nitric oxide synthase III gene in neurodegenerative diseases. *J. Neurol. Sci.* 1999, 162:133-51.
59. Holzmann A, Manktelow C, Taut F, Bloch KD, Zapol WM. Inhibition of nitric oxide synthase prevents hyporesponsiveness to inhaled nitric oxide in lungs from endotoxin-challenged rats. *Anesthesiology* 1999, 91:215-21.
60. Powel V, Moreira GA, O'Donnell DC, Filippov G, Bloch KD, Gordon JB. Maturation changes in ovine pulmonary vascular responses to inhaled nitric oxide. *Pediatric Pulmonology* 1999, 27:157-66.
61. Bloch DB, Chiche J-D, Orth D, Rosenzweig A, Bloch KD. Structural and functional heterogeneity of nuclear bodies. *Molecular and Cellular Biology* 1999, 19:4423-4430.
62. Ullrich R, Bloch KD, Ichinose F, Steudel W, Zapol WM. Preservation of hypoxic pulmonary blood flow redistribution and arterial oxygenation in endotoxin-challenged mice with congenital NOS2 deficiency. *J. Clin. Invest.* 1999, 104:1421-1429.
63. Weimann J, Bloch KD, Takata M, Steudel W, Zapol WM. Congenital NOS2 deficiency protects mice from LPS-induced hyporesponsiveness to inhaled NO. *Anesthesiology* 1999, 91:1744-1753.
64. Weimann J, Ullrich R, Hromi J, Fujino Y, Clark M, Bloch KD, Zapol WM. Sildenafil is a pulmonary vasodilator in awake lambs with acute pulmonary hypertension. *Anesthesiology* 2000, 92:1702-1712.
65. Roberts JD, Chiche J-D, Weimann J, Steudel W, Zapol WM, Bloch KD. Nitric oxide inhalation decreases pulmonary artery remodeling in the injured lungs of rat pups. *Circ. Res.* 2000, 87:140-145.
66. Bloch DB, Nakajima A, Gulick T, Chiche J-D, Orth D, de la Monte SM, Bloch KD. Sp110 localizes to the PML/Sp100 nuclear body and may function as a nuclear hormone receptor transcriptional co-activator. *Molecular and Cellular Biology* 2000, 20:6138-6146.
67. Ullrich R, Scherrer-Crosbie M, Bloch KD, Ichinose F, Nakajima H, Picard MH, Zapol WM, Quezado ZMN. Congenital deficiency of nitric oxide synthase 2 protects against endotoxin-induced myocardial dysfunction in mice. *Circulation* 2000, 102:1440-1446.
68. Budts W, Pokreisz P, Nong Z, Van Pelt N, Gillijns H, Gerard R, Lyons R, Collen D, Bloch KD, Janssens S. Aerosol gene transfer with inducible nitric oxide synthase reduces hypoxic pulmonary hypertension and pulmonary vascular remodeling in rats. *Circulation* 2000, 102:2880-2885.

69. Sinnaeve P, Chiche J-D, Nong Z, Varenne O, Van Pelt N, Gillijns H, Collen D, Bloch KD, Janssens S. Soluble guanylate cyclase  $\alpha 1$  and  $\beta 1$  gene transfer increases NO responsiveness and reduces neointima formation after balloon injury in rats via antiproliferative and antimigratory effects. *Circ. Res.* 2001, 88:103-109.
70. Takata M, Filippov G, Liu H, Ichinose F, Janssens S, Bloch DB, Bloch KD. Cytokines decrease soluble guanylate cyclase in pulmonary artery smooth muscle cells via NO-dependent and NO-independent mechanisms. *Am. J. Physiol.* 2001, 280:L272-L278.
71. Holzmann A, Manktelow C, Weimann J, Bloch KD, Zapol WM. Inhibition of lung phosphodiesterase improves responsiveness to inhaled nitric oxide in isolated-perfused lungs from rats challenged with endotoxin. *Intensive Care Medicine* 2001, 27:251-257.
72. Ichinose F, Zapol WM, Sapirstein A, Ullrich R, Tager AM, Coggins K, Jones R, Bloch KD. Attenuation of hypoxic pulmonary vasoconstriction by endotoxemia requires 5-lipoxygenase in mice. *Circ. Res.* 2001; 88:832-838.
73. Ichinose F, Erana-Garcia J, Hromi J, Raveh Y, Jones R, Krim L, Clark MWH, Winkler JD, Bloch KD, Zapol WM. Nebulized sildenafil is a selective pulmonary vasodilator in lambs with acute pulmonary hypertension. *Critical Care Medicine* 2001, 29:1000-1005.
74. Schmidt U, Han RO, DiSalvo TG, Guerrero JL, Gold HK, Zapol WM, Bloch KD, Semigran MJ. Cessation of platelet-mediated cyclic canine coronary occlusion after thrombolysis by combining nitric oxide inhalation with phosphodiesterase inhibition. *Journal of the American College of Cardiology* 2001, 37:1981-1988.
75. Scherrer-Crosbie M, Ullrich R, Bloch KD, Nakajima H, Aretz HT, Lindsey ML, Vançon A-C, Nasser B, Huang PL, Lee RT, Zapol WM, Picard MH. Nitric oxide synthase 3 limits left ventricular remodeling after myocardial infarction in mice. *Circulation* 2001; 104:1286-1291.
76. Raveh Y, Ichinose F, Orbach P, Bloch KD, Zapol WM. Radical scavengers protect murine lung from endotoxin-induced hyporesponsiveness to inhaled nitric oxide. *Anesthesiology* 2002, 96:926-33.
77. Ichinose F, Ullrich R, Sapirstein A, Jones RC, Bonventre JV, Serhan CM, Bloch KD, Zapol WM. Cytosolic phospholipase A<sub>2</sub> in hypoxic pulmonary vasoconstriction. *J. Clin. Invest.* 2002 109:1493-500.
78. Sinnaeve P, Chiche J-D, Gillijns H, Van Pelt N, Wirthlin DJ, Van de Werf F, Collen D, Bloch KD, Janssens S. Overexpression of a constitutively-active protein kinase G mutant reduces neointima formation and in-stent restenosis. *Circulation* 2002, 105:2911-6.

79. Lepore J, Maroo A, Pereira N, Ginns L, Dec G, Zapol W, Bloch K, Semigran M. Effect of sildenafil on the acute pulmonary vasodilator response to inhaled nitric oxide in adults with primary pulmonary hypertension. *Am J Cardiol* 2002, 90:677-680.
80. Baboolal HA, Ichinose F, Ullrich R, Kawai N, Bloch KD, Zapol WM. Reactive oxygen species scavengers prevent endotoxin-induced impairment of hypoxic pulmonary vasoconstriction in mice. *Anesthesiology* 2002, 97:1227-33.
81. de la Monte SM, Chiche J-D, von dem Bussche A, Sanyal S, Lahousse SA, Janssens SP, Bloch KD. Nitric oxide synthase-3 over-expression causes apoptosis and impairs neuronal mitochondrial function: relevance to Alzheimer-type neurodegeneration. *Lab. Invest.* 2003, 83:287-98.
82. Liu H, Nowak R, Chao W, Bloch KD. Nerve growth factor induces anti-apoptotic heme oxygenase-1 in rat pheochromocytoma PC12 Cells. *J. Neurochem.* 2003, 86:1553-63.
83. Wu JC, Nasser BA, Bloch KD, Picard MH, Scherrer-Crosbie M. Influence of sex on ventricular remodeling after myocardial infarction in mice. *J Am Soc Echocardiogr.* 2003, 16:1158-62.
84. Ichinose F, Hataishi R, Wu JC, Kawai N, Tude-Rodrigues AC, Mallari C, Post JM, Parkinson JF, Picard MH, Bloch KD, Zapol WM. A selective inducible nitric oxide synthase dimerization inhibitor prevents systemic, cardiac and pulmonary hemodynamic dysfunction in endotoxemic mice. *Am J Physiol* 2003, 285:H2524-30.
85. Yu JH, Nakajima A, Nakajima H, Diller LR, Bloch KD, Bloch DB. Restoration of PML-nuclear bodies in neuroblastoma cells enhances retinoic acid responsiveness. *Cancer Res* 2004 64:928-33.
86. Ichinose F, Bloch KD, Wu JC, Hataishi R, Aretz HT, Picard MH, Scherrer-Crosbie M. Pressure-overload induced left ventricular hypertrophy and dysfunction in mice are exacerbated by congenital nitric oxide synthase 3 deficiency. *Am J Physiol* 2004 286:H1070-5.
87. Hassoun PM, Filippov G, Fogel M, Donaldson C, Kayyali US, Shimoda LA, Bloch KD. Hypoxia decreases expression of soluble guanylate cyclase in cultures rat pulmonary artery smooth muscle cells. *Am J Respir Cell Mol Biol.* 2004, 30:908-13.
88. Janssens S, Pokreisz P, Schoonjans L, Pellens M, Vermeersch P, Tjwa M, Jans P, Scherrer-Crosbie M, Picard MH, Szelis Z, Gillijns H, Van de Werf F, Collen D, Bloch KD. Cardiomyocyte-specific over-expression of nitric oxide synthase 3 improves left ventricular performance and reduces compensatory hypertrophy after myocardial infarction. *Circ Res*, 2004, 94:1256-62.

89. del Monte F, Dalal R, Tabchy A, Couget J, Bloch KD, Peterson R, Hajjar RJ. Transcriptional changes following restoration of SERCA2a levels in failing rat hearts. *FASEB J.* 2004 18:1474-6.
90. Inglessis I, Shin JT, Lepore JJ, Palacios IF, Zapol WM, Bloch KD, Semigran MJ. Hemodynamic effects of inhaled nitric oxide in right ventricular myocardial infarction and cardiogenic shock. *Journal of the American College of Cardiology* 2004 44:793-8.
91. Tudes-Rodrigues AC, Hataishi R, Ichinose F, Bloch KD, Derumeaux G, Picard MH, Scherrer-Crosbie M. Relationship of systolic dysfunction to area at risk and infarction size after ischemia-reperfusion in mice. *J Am Soc Echocardiogr* 2004, 17:948-53.
92. Lewis G, Bloch KD, Semigran MJ. Pulmonary thromboembolism superimposed on a congenital VSD in a 50 year old man; inhaled nitric oxide and sildenafil to the rescue. *Cardiology in Review* 2004, 12:188-90.
93. Beppu H, Ichinose F, Kawai N, Jones RC, Yu P, Zapol WM, Miyazono K, Li E, Bloch KD. BMPR-II heterozygous mice have mild pulmonary hypertension and an impaired pulmonary vascular remodeling response to prolonged hypoxia. *Am J Physiol* 2004, 287:L1241-7.
94. Browner NC, Dey NB, Bloch KD, Lincoln TM. Regulation of cGMP-dependent protein kinase expression by soluble guanylyl cyclase in vascular smooth muscle cells. *J Biol Chem.* 2004, 279:46631-6.
95. Evgenov OV, Ichinose F, Evgenov NV, Gnoth MJ, Falkowski GE, Chang Y, Bloch KD, Zapol WM. Soluble guanylate cyclase activator reverses acute pulmonary hypertension and augments the pulmonary vasodilator response to inhaled nitric oxide in awake lambs. *Circulation* 2004, 110:2253-9.
96. Bloch DB, Yu JH, Yang W-H, Graeme-Cook F, Lindor K, Viswanathan A, Bloch KD, Nakajima A. The cytoplasmic dot staining pattern is detected in a subgroup of patients with primary biliary cirrhosis. *Journal of Rheumatology*, 2005, 32:477-83.
97. Weinberg EO, Scherrer-Crosbie M, Picard MH, Nasser BA, MacGillivray C, Gannon J, Lian Q, Bloch KD, Lee RT. Rosuvastatin reduces experimental left ventricular infarct size following ischemia-reperfusion injury but not total coronary occlusion. *Am J Physiol Heart Circ Physiol.* 2005, 288:H1802-9.
98. Beppu H, Lei H, Bloch KD, Li E. Generation of a floxed allele of the mouse BMP type II receptor gene. *Genesis* 2005, 41:133-7.
99. Lepore JJ, Maroo A, Pereira NL, Bigatello LM, Dec GW, Zapol WM, Bloch KD, Semigran MJ. Hemodynamic effects of sildenafil in patients with congestive heart failure and pulmonary hypertension: combined administration with inhaled nitric oxide. *Chest*, 2005, 127:1647-53.

100. Sebag IA, Handschumacher MD, Ichinose F, Morgan JG, Hataishi R, Rodrigues ACT, Guerrero JL, Steudel W, Raher MJ, Halpern EF, Derumeaux G, Bloch KD, Picard MH, and Scherrer-Crosbie M. Quantitative assessment of regional myocardial function in mice by tissue doppler imaging: comparison with hemodynamics and sonomicrometry. *Circulation* 2005, 111:2611-6.
101. Yu PB, Beppu H, Kawai N, Li E, Bloch KD. BMP type II receptor deletion reveals BMP ligand-specific gain of signaling in pulmonary artery smooth muscle cells. *J Biol Chem* 2005, 280:24443-24450.
102. Lepore JJ, Dec GW, Zapol WM, Bloch KD, Semigran MJ. Combined administration of intravenous dipyridamole and inhaled nitric oxide to assess reversibility of pulmonary arterial hypertension in potential transplant patients. *Journal of Heart and Lung Transplantation* 2005, in press.
103. Caironi P, Ichinose F, Liu R, Jones RC, Bloch KD, Zapol WM. 5-lipoxygenase deficiency prevents respiratory failure during ventilator-induced lung injury. *American Journal of Respiratory and Critical Care Medicine* 2005, in press.

#### **Proceedings of Meetings:**

1. Bloch KD, Seidman JG, Seidman CE. Structure and expression of the atrial natriuretic factor gene. In: *Atrial hormones and other natriuretic factors. Clinical Physiology Series, American Physiological Society, Bethesda, MD, 1987, p. 7-18.*

#### **Reviews, Chapters, and Editorials:**

1. Seidman CE, Bloch KD, Zisfein JB, Smith JA, Haber E, Homcy CJ, Duby AD, Choi E, Graham RM, Seidman JG. Molecular studies of the atrial natriuretic factor gene. *Hypertension* 1985; 7 (part II): 31-34.
2. Seidman CE, Bloch KD. Molecular approaches to the study of atrial natriuretic factor. *Am J Med Sci* 1987; 294:144-149.
3. Neer EJ, Kim SY, Ang SL, Bloch DB, Bloch KD, Kawahara K, Tolman C, Lee R, Logethetis D, Kim D, Seidman JG, Clapham DE. Functions of G protein subunits. *Cold Spring Harbor Symp. Quant. Biol.* 1989; 53:241-246.
4. Bloch KD. The nitric oxide synthase gene family. In: *Molecular Biology of the Kidney in Health and Disease. Schlondorff D and Bonventre JV, ed., Marcel Dekker Inc. New York, NY, chapter 11, pp 133-142, 1995.*
5. Lloyd-Jones D, Bloch KD. Vascular biology of nitric oxide and its role in atherogenesis. In: *Annual Review of Medicine. C. Coggins, editor. Annual Reviews Inc. Palo Alto CA, volume 47, p. 365-376, 1996.*

6. Lepore JJ, Bloch KD. Nitric oxide and pulmonary hypertension. In: Loscalzo J and Vita JA, eds. Nitric oxide and the cardiovascular system. Totowa, NJ: Humana Press, Inc., 1999, pp. 247-272.

7. Bloch KD. Regulation of endothelial NO synthase mRNA stability: RNA-binding proteins crowd on the 3'-untranslated region. Circ. Res. 1999, 85:653-655.

8. Passeri J, Bloch, KD. Nitric Oxide and Cardiac Remodeling. In: Hajjar R, Del Monte F, eds. Heart Failure Clinics. Elsevier Science, Inc. 2005, in press.

**Books, Monographs, and Text Books:**

1. Zapol WM, Bloch KD, editors. Nitric Oxide and the Lung for the series Lung Biology in Health and Disease. Lenfant C, executive editor. Marcel Dekker Inc. New York, NY, 1996.

**Clinical Communications:** none

**Educational Material:** none

**Thesis:** none

**Nonprint Materials:** none

**Patents:**

- 5,904,938: Treatment of vascular thrombosis and restenosis with inhaled nitric oxide.
- 6,063,407: Treatment of vascular thrombosis and restenosis with inhaled nitric oxide.
- 6,183,988: Leukocyte-specific protein and gene, and methods of use thereof.
- 6,601,580: Enhancing therapeutic effectiveness of nitric oxide inhalation.
- 6,656,452: Use of inhaled NO as anti-inflammatory agent.
- 6,720,309: Method of inducing vasodilation and treating pulmonary hypertension using adenoviral-mediated transfer of the nitric oxide synthase gene.
- 6,811,768: Use of inhaled NO as anti-inflammatory agent.

# EXHIBIT B



## EXHIBIT B

## FIGURE 1

Alignment of human DNMT3A from ATCC Deposit No. 98809 (top) and currently amended SEQ ID NO:3 (bottom)<sup>1</sup>

```

gccgcggcaccagggcgcgagccgggcccggcccgacccaccggccatacgggtggagcc
|||||
gccgcggcaccagggcgcgagccgggcccggcccgacccaccggccatacgggtggagcc 60

atcgaagccccacccacaggctgacagaggcaccgttcaccagaggggtcaacaccggg
|||||
atcgaagccccacccacaggctgacagaggcaccgttcaccagaggggtcaacaccggg 120

atctatgtttaagttttaactctcgctccaaagaccacgataattccttccccaaagcc
|||||
atctatgtttaagttttaactctcgctccaaagaccacgataattccttccccaaagcc 180

cagcagccccccagccccgcgcagccccagcctgcctcccggcgcccagatgcccgccat
|||||
cagcagccccccagccccgcgcagccccagcctgcctcccggcgcccagatgcccgccat 240

gccctccagcggccccggggacaccagcagctctgctgaggagcgggaggaggaccgaaa
|||||
gccctccagcggccccggggacaccagcagctctgctgaggagcgggaggaggaccgaaa 300

ggacggagaggagcaggaggagccgcgtggcaaggaggagcgccaagagcccagcaccac
|||||
ggacggagaggagcaggaggagccgcgtggcaaggaggagcgccaagagcccagcaccac 360

ggcacggaaggtggggcggcctgggaggaagcgcaagcacccccgggtggaaagcgggtga
|||||
ggcacggaaggtggggcggcctgggaggaagcgcaagcacccccgggtggaaagcgggtga 420

cacgccaaaggaccctgcggtgatctccaagtccccatccatggcccaggactcaggcgc
|||||
cacgccaaaggaccctgcggtgatctccaagtccccatccatggcccaggactcaggcgc 480

ctcagagctattacccaatggggacttggagaagcggagttagccccagccagaggag...
|||||
ctcagagctattacccaatggggacttggagaagcggagttagccccagccagaggagg 540

.....agggtgcagctgagac
|||||
gagccctgctggggggcagaagggcggggccccagcagaggagagggtgcagctgagac 600

cctgcctgaagcctcaagagcagtggaagaaatggctgctgcaccccccaaggagggccgagg
|||||
cctgcctgaagcctcaagagcagtggaagaaatggctgctgcaccccccaaggagggccgagg 660

agcccctgcagaagcggggcaaagaacagaaggagaccaacatcgaatccatgaaaatgga
|||||
agcccctgcagaagcggggcaaagaacagaaggagaccaacatcgaatccatgaaaatgga 720

```

<sup>1</sup> Bolded nucleotides indicate nucleotides that were amended on July 23, 2001.

gggctcccggggccgggtgcgggggtggcttgggctgggagtcacagcctccgtcagcggcc  
gggctcccggggccgggtgcgggggtggcttgggctgggagtcacagcctccgtcagcggcc 780

catgccgaggctcaccttccaggcgggggacccctactacatcagcaagcgcaagcggga  
catgccgaggctcaccttccaggcgggggacccctactacatcagcaagcgcaagcggga 840

cgagtggctggcacgctggaaaagggaggctgagaagaaagccaaggtcattgcaggaat  
cgagtggctggcacgctggaaaagggaggctgagaagaaagccaaggtcattgcaggaat 900

gaatgctgtggaagaaaaccagggggccggggagtgctcagaaggtggaggaggccagccc  
gaatgctgtggaagaaaaccagggggccggggagtgctcagaaggtggaggaggccagccc 960

tcttgctgtgcagcagcccactgaccccgcatccccactgtggctaccacgcctgagcc  
tcttgctgtgcagcagcccactgaccccgcatccccactgtggctaccacgcctgagcc 1020

cgtaggggtccgatgctggggacaagaatgccaccaaagcaggcgatgacgagccagagta  
cgtaggggtccgatgctggggacaagaatgccaccaaagcaggcgatgacgagccagagta 1080

cgaggacggccggggccttggcattggggagctggtgtgggggaaactgcggggcttctc  
cgaggacggccggggccttggcattggggagctggtgtgggggaaactgcggggcttctc 1140

ctggtggccaggccgcattgtgtcttgggtggatgacgggcccggagccgagcagctgaagg  
ctggtggccaggccgcattgtgtcttgggtggatgacgggcccggagccgagcagctgaagg 1200

cacccgctgggtcatgtggttcggagacggcaaattctcagtgggtgtgtgttgagaagct  
cacccgctgggtcatgtggttcggagacggcaaattctcagtgggtgtgtgttgagaagct 1260

gatgccgctgagctcgttttgcagtgcgttcaccaggccacgtacaacaagcagcccat  
gatgccgctgagctcgttttgcagtgcgttcaccaggccacgtacaacaagcagcccat 1320

gtaccgcaaagccatctacgaggtcctgcaggtggccagcagccgcgcgggggaagctggt  
gtaccgcaaagccatctacgaggtcctgcaggtggccagcagccgcgcgggggaagctggt 1380

cccgggtgtgccacgacagcgatgagagtgaactgccaaggccgtggaggtgcagaacaa  
cccgggtgtgccacgacagcgatgagagtgaactgccaaggccgtggaggtgcagaacaa 1440

gcccattgattgaatgggccctggggggccttcagccttctggccctaagggcctggagcc  
gcccattgattgaatgggccctggggggccttcagccttctggccctaagggcctggagcc 1500

accagaagaagagaagaatccctacaaagaagtgtacacggacatgtgggtggaacctga  
|||||  
accagaagaagagaagaatccctacaaagaagtgtacacggacatgtgggtggaacctga 1560

ggcagctgcctacgcaccacctccaccagccaaaaagccccggaagagcacagcggagaa  
|||||  
ggcagctgcctacgcaccacctccaccagccaaaaagccccggaagagcacagcggagaa 1620

gccaaggtcaaggagattattgatgagcgcacaagagagcggctggtgtacgaggtgcg  
|||||  
gccaaggtcaaggagattattgatgagcgcacaagagagcggctggtgtacgaggtgcg 1680

gcagaagtgccggaacattgaggacatctgcatctcctgtgggagcctcaatgttaccct  
|||||  
gcagaagtgccggaacattgaggacatctgcatctcctgtgggagcctcaatgttaccct 1740

ggaacacccccctcttcgttggaggaatgtgccaaaactgcaagaactgctttctggagtg  
|||||  
ggaacacccccctcttcgttggaggaatgtgccaaaactgcaagaactgctttctggagtg 1800

tgcgtagcagtagcagcagcagcggctaccagtcctactgcaccatctgctgtggggggcgg  
|||||  
tgcgtagcagtagcagcagcagcggctaccagtcctactgcaccatctgctgtggggggcgg 1860

tgaggtgctcatgtgcggaacaacaactgctgcaggtgcttttgctggagtggtgga  
|||||  
tgaggtgctcatgtgcggaacaacaactgctgcaggtgcttttgctggagtggtgga 1920

cctcttggtggggccgggggctgccaggcagccattaaggaagacccctggaactgcta  
|||||  
cctcttggtggggccgggggctgccaggcagccattaaggaagacccctggaactgcta 1980

catgtgcgggcacaaaggtacctacgggctgctgcggcggcgagaggactggccctcccg  
|||||  
catgtgcgggcacaaaggtacctacgggctgctgcggcggcgagaggactggccctcccg 2040

gctccagatgttcttcgctaataaccacgaccaggaatttgacctccaaaggtttacc  
|||||  
gctccagatgttcttcgctaataaccacgaccaggaatttgacctccaaaggtttacc 2100

acctgtcccagctgagaagaggaagcccatccgggtgctgtctctctttgatggaatcgc  
|||||  
acctgtcccagctgagaagaggaagcccatccgggtgctgtctctctctttgatggaatcgc 2160

tacagggctcctggtgctgaaggacttgggcattcaggtggaccgctacattgcctcgga  
|||||  
tacagggctcctggtgctgaaggacttgggcattcaggtggaccgctacattgcctcgga 2220

ggtgtgtgaggactccatcacggtgggcatggtgcggcaccaggggaagatcatgtacgt  
|||||  
ggtgtgtgaggactccatcacggtgggcatggtgcggcaccaggggaagatcatgtacgt 2280

cggggacgtccgcagcgtcacacagaagcatatccaggagtggggcccatcgcattctggt  
|||||  
cggggacgtccgcagcgtcacacagaagcatatccaggagtggggcccatcgcattctggt 2340

gattgggggcagtccttgcaatgacctctccatcgtaaccctgctcgcaagggcctcta  
|||||  
gattgggggcagtccttgcaatgacctctccatcgtaaccctgctcgcaagggcctcta 2400

cgagggcactggccggctcttctttgagttctaccgcctcctgcatgatgcgcgcccaa  
|||  
cgagggcactggccggctcttctttgagttctaccgcctcctgcatgatgcgcgcccaa 2460

ggagggagatgatcgccccttcttctggctctttgagaatgtggtggccatgggcgtag  
|||  
ggagggagatgatcgccccttcttctggctctttgagaatgtggtggccatgggcgtag 2520

tgacaagagggacatctcgcgatttctcgagtccaaccctgtgatgattgatgccaaaga  
|||  
tgacaagagggacatctcgcgatttctcgagtccaaccctgtgatgattgatgccaaaga 2580

agtgtcagctgcacacagggcccgtacttctggggtaaccttcccggtatgaacaggcc  
|||  
agtgtcagctgcacacagggcccgtacttctggggtaaccttcccggtatgaacaggcc 2640

gttggcatccactgtgaatgataagctggagctgcaggagtgctctggagcatggcaggat  
|||  
gttggcatccactgtgaatgataagctggagctgcaggagtgctctggagcatggcaggat 2700

agccaagttcagcaaagtgaggaccattactacgaggtcaaactccataaagcagggcaa  
|||  
agccaagttcagcaaagtgaggaccattactacgaggtcaaactccataaagcagggcaa 2760

agaccagcatttttctgtcttcatgaatgagaaagaggacatcttatggtgcaactgaaat  
|||  
agaccagcatttttctgtcttcatgaatgagaaagaggacatcttatggtgcaactgaaat 2820

ggaaagggatatttggtttccagtcactatactgacgtctccaacatgagccgcttggc  
|||  
ggaaagggatatttggtttccagtcactatactgacgtctccaacatgagccgcttggc 2880

gagggcagagactgctgggcccgtcatggagcgtgccagtcacccgccacctcttcgctcc  
|||  
gagggcagagactgctgggcccgtcatggagcgtgccagtcacccgccacctcttcgctcc 2940

gctgaaggagtatatttgctgtgtgtgaagggacatgggggcaaactgaggtagcg  
|||  
gctgaaggagtatatttgctgtgtgtgaagggacatgggggcaaactgaggtagcg 2995

# FIGURE 2

Alignment of predicted amino acids encoded by DNMT3A cDNA in ATCC Deposit No. 98809 (top) and predicted amino acids encoded by currently amended SEQ ID NO:3 (bottom)<sup>2</sup>

MPAMPSSGPGDTSSSSAEREEDRKDGEEQEEPRGKEERQEPSTTARKVGRPGKR	55
MPAMPSSGPGDTSSSSAEREEDRKDGEEQEEPRGKEERQEPSTTARKVGRPGKR	55
KHPPVESGDTPKDPAVISKSPSMAQDSGASELLPNGDLEKRSEPQPEERVQLRPC	110
KHPPVESGDTPKDPAVISKSPSMAQDSGASELLPNGDLEKRSEPQPEEGSPAGGQ	110
<b>LKPQEQWKMAAAPPRRAEEPLQKRAKNRRRPTSNP*KWRAPGAGCGVAWAGSPAS</b>	165
KGGAPAEEGEAAETLPEASRAVENGCCTPKEGRGAPAEAGKEQKETNIESMKMEG	165
VSGPCRGSPPRRGTPTTSASASGTSQWHAGKGRLLRRKPRSLQE*MLWKKTRGPGS	220
SRGRLRGGLGWESSLRQRPMPRLTFQAGDPYYISKRRDEWLARWKREAEKKAKV	220
LRRWRRPALLLCSSPLTPHPPLWLPRLSWPWGPMLGTRMPPKQAMTSQSSTRTAGAL	275
GMNAVEENQGPGESQKVEEASPPAVQQPTDPASPTVATTPEPVGSDAGDKNIAAT	275
ALGSWCGGNCGASPGGQAALCLGG*AGAEQLKAPAGSCGSETANSQWCVLRS*C	330
KAGDDEPEYEDGRGFGIGELVWGKLRGFSWWPGRIVSWWMTGRSRAAEGTRWVMW	330
<b>R*ARFAVRSTRPRTTSSPCTAKPSTRSCRWPAARAGSCSRCATTAMRVTLPRPWR</b>	385
FGDGKFSVVCVEKLMPLSSFCFAHQATYNKQPMYRKAIYEVQLQVASSRAGKLF	385
CRTSP*LNGPWGASSLLALRAWSHQKKRIPTKKCTRTCGWNLRQLPTHHLHQPK	440
VCHDSDESDTAKAVEVQNKPMIEWALGGFQPSGPKGLEPPEEEKNPYKEVYTDW	440
SPGRAQRSPRSRRLMSAQESGWCTRCGRSAGTLRTSASPVGASMLPWNTPSLL	495
VEPEAAAYAPPPPAKKPRKSTAEKPKVKEIIDERTRELRVYEVQRKCRNIEDICI	495
ECAKTARTAFWSVRTSTTTTATSPTAPSAVGAVRCSAETTTAAGAFASVWTS	550
SCGSLNVTLEHPLFVGGMCQNCNCFLECAQYQDDGYQSYCTICCGGREVLNMG	550
WWGRGLPRQPLRKTPGTATCAGTRVPTGCCGGERTGPPGSRCSLLITTRNLTLQ	605
NNNCCRCFCVECDLLVGPAAQAAIKEDPWNCYMCCHKGTGGLRRREDWPSRL	605
RFTHLSQLRRGSPSGCCLSLMESLQGSWC*RTWAFRWTATLPRRCVTRTPSRWAWC	660
QMFFANNHDQEFDPKVPVPAEKRPPIRVLSLFDGIATGLLVKDLGIQVDY	660
GTRGRSCTSGTSAASHRSISRGAHSIW*LGAVPAMTSPSSTLLARASTRALAGS	715
IASEVCEDSITVGMVRHQKIMYVGDVRSVTQKHIQEWGPFDLVIGGSPCNDLSI	715
SLSTASCMMRGPRREMIAPSSGSLRMWWPWALVTRGTSRDFSSPTL**LMPKKC	770
VNPARKGLYEGTGRLFFEFYRLLDARPKEGDDRPFFWLFENVVAMGVSDKRDIS	770
QLHTGPATSGVTFFV*TGRWHPL*MISWSCRSVWSMAG*PSSAK*GPLLRGQTP*	825
RFLESNPVMIDAKEVSAAHRARYFWGNLPGMNRPLASTVNDKLELQECLEHGRIA	825

<sup>2</sup> \* indicates a predicted stop codon. Bolded amino acids are encoded by nucleotides located downstream of the deletion.

SRAKTSIFLSS*MRKRTSYGALKWKGYLVSQSTILTSPT*AAWRGRDCWAGHGAC	880
KFSKVRTITTRSNSIKQKQDQHFPVFMNEKEDILWCTEMERVFGFPVHYTDVSNM	880
QSSATSSLR*RSILRVCKGHGGKLR**AAWRG	912
SRLARQRLGRSWSVPVIRHLFAPLKEYFACV	912

FIGURE 3

Human DNMT3A from ATCC Deposit No. 98809 with forward reading frames  
(DNMT3A amino acid residues bolded)

DNA: GCCGCGGCACCAGGGCGCGCAGCCGGGCGGCCGACCCACCGGCCATAC 51  
+3: R G T R A R S R A G P T P P A I R  
+2: P R H Q G A Q P G R P D P T G H T  
+1: A A A P G R A A G P A R P H R P Y

DNA: GGTGGAGCCATCGAAGCCCCACCCACAGGCTGACAGAGGCACCGTTTACC 102  
+3: W S H R S P H P Q A D R G T V H Q  
+2: V E P S K P P P T G \* Q R H R S P  
+1: G G A I E A P T H R L T E A P F T

DNA: AGAGGGCTCAACACCGGGATCTATGTTTAAAGTTTAACTCTCGCCTCCAAA 153  
+3: R A Q H R D L C L S F N S R L Q R  
+2: E G S T P G S M F K F \* L S P P K  
+1: R G L N T G I Y V \* V L T L A S K

DNA: GACCACGATAATTCCTTCCCCAAAGCCCAGCAGCCCCCAGCCCCGCGCAG 204  
+3: P R \* F L P Q S P A A P Q P R A A  
+2: T T I I P S P K P S S P P A P R S  
+1: D H D N S F P K A Q Q P P S P A Q

DNA: CCCCAGCTGCCTCCCGGCGCCCATGCCCCGCGCCATGCCCTCCAGCGGCC 255  
+3: P A C L P A P R C P P C P P A A P  
+2: P S L P P G A Q M P A M P S S G P  
+1: P Q P A S R R P D A R H A L Q R P

DNA: CGGGGACACCAGCAGCTCTGCTGCGGAGCGGGAGGAGGACCGAAAGGACGG 306  
+3: G T P A A L L R S G R R T E R T E  
+2: G D T S S S A A E R E E D R K D G  
+1: R G H Q Q L C C G A G G G P K G R

DNA: AGAGGAGCAGGAGGAGCCGCGTGGCAAGGAGGAGCGCCAAGAGCCCAGCAC 357  
+3: R S R R S R V A R R S A K S P A P  
+2: E E Q E E P R G K E E R Q E P S T  
+1: R G A G G A A W Q G G A P R A Q H

DNA: CACGGCACGGAAGGTGGGGCGGCCTGGGAGGAAGCGCAAGCACCCCCCGGT 408  
+3: R H G R W G G L G G S A S T P R W  
+2: T A R K V G R P G R K R K H P P V  
+1: H G T E G G A A W E E A Q A P P G

DNA: GGAAAGCGGTGACACGCCAAAGGACCCTGCGGTGATCTCCAAGTCCCCATC 459  
+3: K A V T R Q R T L R \* S P S P H P  
+2: E S G D T P K D P A V I S K S P S  
+1: G K R \* H A K G P C G D L Q V P I

DNA: CATGGCCCAGGACTCAGGCGCCTCAGAGCTATTACCCAATGGGGACTTGGA 510  
+3: W P R T Q A P Q S Y Y P M G T W R  
+2: M A Q D S G A S E L L P N G D L E  
+1: H G P G L R R L R A I T Q W G L G

DNA: GAAGCGGAGTGAGCCCCAGCCAGAGGAG..... 561  
+3: S G V S P S Q R R.....  
+2: K R S E P Q P E E.....DELETION.....  
+1: E A E \* A P A R G E.....

DNA: .....AGGGTGCAGCTGAGACCCCTGCCTGAAGC 612  
+3: .....G C S \* D P A \* S  
+2: .....R V Q L R P C L K P  
+1: .....G A A E T L P E A

DNA: CTCAAGAGCAGTGGAAAATGGCTGCTGCACCCCAAGGAGGGCCGAGGAGC 663  
+3: L K S S G K W L L H P Q G G P R S  
+2: Q E Q W K M A A A P P R R A E E P  
+1: S R A V E N G C C T P K E G R G A

DNA: CCCTGCAGAAGCGGGCAAAGAACAGAAGGAGACCAACATCGAATCCATGAA 714  
+3: P C R S G Q R T E G D Q H R I H E  
+2: L Q K R A K N R R R P T S N P \* K  
+1: P A E A G K E Q K E T N I E S M K

DNA: AATGGAGGGCTCCCGGGGCGGCTGCGGGGTGGCTTGGGCTGGGAGTCCAG 765  
+3: N G G L P G P A A G W L G L G V Q  
+2: W R A P G A G C G V A W A G S P A  
+1: M E G S R G R L R G G L G W E S S

DNA: CCTCCGTCAGCGGCCCATGCCGAGGCTCACCTTCCAGGCGGGGACCCCTA 816  
+3: P P S A A H A E A H L P G G G P L  
+2: S V S G P C R G S P S R R G T P T  
+1: L R Q R P M P R L T F Q A G D P Y

DNA: CTACATCAGCAAGCGCAAGCGGGACGAGTGGCTGGCACGCTGGAAGGGA 867  
+3: L H Q Q A Q A G R V A G T L E K G  
+2: T S A S A S G T S G W H A G K G R  
+1: Y I S K R K R D E W L A R W K R E

DNA: GGCTGAGAAGAAAGCCAAGGTCATTGCAGGAATGAATGCTGTGGAAGAAAA 918  
+3: G \* E E S Q G H C R N E C C G R K  
+2: L R R K P R S L Q E \* M L W K K T  
+1: A E K K A K V I A G M N A V E E N

DNA: CCAGGGGCGCGGGAGTCTCAGAAGGTGGAGGAGGCCAGCCCTCCTGCTGT 969  
+3: P G A R G V S E G G G G Q P S C C  
+2: R G P G S L R R W R R P A L L L C  
+1: Q G P G E S Q K V E E A S P P A V

DNA: GCAGCAGCCCACTGACCCCGCATCCCCACTGTGGCTACCACGCCTGAGCC 1020  
+3: A A A H \* P R I P H C G Y H A \* A  
+2: S S P L T P H P P L W L P R L S P  
+1: Q Q P T D P A S P T V A T T P E P

DNA: CGTGGGGTCCGATGCTGGGGACAAGAATGCCACCAAAGCAGGCGATGACGA 1071  
+3: R G V R C W G Q E C H Q S R R \* R  
+2: W G P M L G T R M P P K Q A M T S  
+1: V G S D A G D K N A T K A G D D E

DNA: GCCAGAGTACGAGGACGGCCGGGGCTTTGGCATTGGGGAGCTGGTGTGGGG 1122  
+3: A R V R G R P G L W H W G A G V G  
+2: Q S T R T A G A L A L G S W C G G  
+1: P E Y E D G R G F G I G E L V W G

DNA: GAAACTGCGGGGCTTCTCCTGGTGGCCAGGCCGATTGTGTCTTGGTGGAT 1173  
+3: E T A G L L L V A R P H C V L V D  
+2: N C G A S P G G Q A A L C L G G \*  
+1: K L R G F S W W P G R I V S W W M



DNA: GACGGGCCGAGCCGAGCAGCTGAAGGCACCCGCTGGGTCATGTGGTTTCGG 1224  
+3: D G P E P S S \* R H P L G H V V R  
+2: R A G A E Q L K A P A G S C G S E  
+1: T G R S R A A E G T R W V M W F G

DNA: AGACGGCAAATTCTCAGTGGTGTGTGTTGAGAAGCTGATGCCGCTGAGCTC 1275  
+3: R R Q I L S G V C \* E A D A A E L  
+2: T A N S Q W C V L R S \* C R \* A R  
+1: D G K F S V V C V E K L M P L S S

DNA: GTTTTGCAGTGC GTTCCACCAGGCCACGTACAACAAGCAGCCCATGTACCG 1326  
+3: V L Q C V P P G H V Q Q A A H V P  
+2: F A V R S T R P R T T S S P C T A  
+1: F C S A F H Q A T Y N K Q P M Y R

DNA: CAAAGCCATCTACGAGGTCTGTCAGGTGGCCAGCAGCCGCGGGGAAGCT 1377  
+3: Q S H L R G P A G G Q Q P R G E A  
+2: K P S T R S C R W P A A A R G S C  
+1: K A I Y E V L Q V A S S R A G K L

DNA: GTTCCCGGTGTGCCACGACAGCGATGAGAGTGACACTGCCAAGGCCGTGGA 1428  
+3: V P G V P R Q R \* E \* H C Q G R G  
+2: S R C A T T A M R V T L P R P W R  
+1: F P V C H D S D E S D T A K A V E

DNA: GGTGCAGAACAAAGCCCATGATTGAATGGGCCCTGGGGGGCTTCCAGCCTTC 1479  
+3: G A E Q A H D \* M G P G G L P A F  
+2: C R T S P \* L N G P W G A S S L L  
+1: V Q N K P M I E W A L G G F Q P S

DNA: TGGCCCTAAGGGCCTGGAGCCACCAGAAGAAGAGAAGAATCCCTACAAAGA 1530  
+3: W P \* G P G A T R R R E E S L Q R  
+2: A L R A W S H Q K K R R I P T K K  
+1: G P K G L E P P E E E K N P Y K E

DNA: AGTGTACACGGACATGTGGGTGGAACCTGAGGCAGCTGCCTACGCACCACC 1581  
+3: S V H G H V G G T \* G S C L R T T  
+2: C T R T C G W N L R Q L P T H H L  
+1: V Y T D M W V E P E A A A Y A P P

DNA: TCCACCAGCCAAAAAGCCCCGGAAGAGCACAGCGGAGAAGCCCAAGGTCAA 1632  
+3: S T S Q K A P E E H S G E A Q G Q  
+2: H Q P K S P G R A Q R R S P R S R  
+1: P P A K K P R K S T A E K P K V K

DNA: GGAGATTATTGATGAGCGCACAAAGAGAGCGGCTGGTGTACGAGGTGCGGCA 1683  
+3: G D Y \* \* A H K R A A G V R G A A  
+2: R L L M S A Q E S G W C T R C G R  
+1: E I I D E R T R E R L V Y E V R Q

DNA: GAAGTGCCGGAACATTGAGGACATCTGCATCTCCTGTGGGAGCCTCAATGT 1734  
+3: E V P E H \* G H L H L L W E P Q C  
+2: S A G T L R T S A S P V G A S M L  
+1: K C R N I E D I C I S C G S L N V

DNA: TACCCTGGAACACCCCTCTTCGTTGGAGGAATGTGCCAAAAGTCAAGAA 1785  
+3: Y P G T P P L R W R N V P K L Q E  
+2: P W N T P S S L E E C A K T A R T  
+1: T L E H P L F V G G M C Q N C K N

DNA: CTGCTTTCTGGAGTGTGCGTACCAGTACGACGACGACGGCTACCAGTCCTA 1836  
+3: L L S G V C V P V R R R R L P V L  
+2: A F W S V R T S T T T T A T S P T  
+1: C F L E C A Y Q Y D D D G Y Q S Y

DNA: CTGCACCATCTGCTGTGGGGGCCGTGAGGTGCTCATGTGCGGAAACAACAA 1887  
+3: L H H L L W G P \* G A H V R K Q Q  
+2: A P S A V G A V R C S C A E T T T  
+1: C T I C C G G R E V L M C G N N N

DNA: CTGCTGCAGGTGCTTTTTCGTGGAGTGTGTGGACCTCTTGGTGGGGCCGGG 1938  
+3: L L Q V L L R G V C G P L G G A G  
+2: A A G A F A W S V W T S W W G R G  
+1: C C R C F C V E C V D L L V G P G

DNA: GGCTGCCCAGGCAGCCATTAAGGAAGACCCCTGGAAGTCTACATGTGCGG 1989  
+3: G C P G S H \* G R P L E L L H V R  
+2: L P R Q P L R K T P G T A T C A G  
+1: A A Q A A I K E D P W N C Y M C G

DNA: GCACAAGGGTACCTACGGGCTGCTGCGGCGGCGAGAGGACTGGCCCTCCCG 2040  
+3: A Q G Y L R A A A A A R G L A L P  
+2: T R V P T G C C G G E R T G P P G  
+1: H K G T Y G L L R R R E D W P S R

DNA: GCTCCAGATGTTCTTCGCTAATAACCACGACCAGGAATTTGACCCTCCAAA 2091  
+3: A P D V L R \* \* P R P G I \* P S K  
+2: S R C S S L I T T T R N L T L Q R  
+1: L Q M F F A N N H D Q E F D P P K

DNA: GGTTTACCCACCTGTCCCAGCTGAGAAGAGGAAGCCCATCCGGGTGCTGTC 2142  
+3: G L P T C P S \* E E E A H P G A V  
+2: F T H L S Q L R R G S P S G C C L  
+1: V Y P P V P A E K R K P I R V L S

DNA: TCTCTTTGATGGAATCGCTACAGGGCTCCTGGTGTGTAAGGACTTGGGCAT 2193  
+3: S L \* W N R Y R A P G A E G L G H  
+2: S L M E S L Q G S W C \* R T W A F  
+1: L F D G I A T G L L V L K D L G I

DNA: TCAGGTGGACCGCTACATTGCCTCGGAGGTGTGTGAGGACTCCATCACGGT 2244  
+3: S G G P L H C L G G V \* G L H H G  
+2: R W T A T L P R R C V R T P S R W  
+1: Q V D R Y I A S E V C E D S I T V

DNA: GGGCATGGTGCGGCACCAGGGGAAGATCATGTACGTGCGGGACGTCCGCAG 2295  
+3: G H G A A P G E D H V R R G R P Q  
+2: A W C G T R G R S C T S G T S A A  
+1: G M V R H Q G K I M Y V G D V R S

DNA: CGTCACACAGAAGCATATCCAGGAGTGGGGCCCATTTCGATCTGGTGATTGG 2346  
+3: R H T E A Y P G V G P I R S G D W  
+2: S H R S I S R S G A H S I W \* L G  
+1: V T Q K H I Q E W G P F D L V I G

DNA: GGGCAGTCCCTGCAATGACCTCTCCATCGTCAACCCTGCTCGCAAGGGCCT 2397  
+3: G Q S L Q \* P L H R Q P C S Q G P  
+2: A V P A M T S P S S T L L A R A S  
+1: G S P C N D L S I V N P A R K G L

DNA: CTACGAGGGCACTGGCCGGCTCTTCTTTGAGTTCTACCGCCTCCTGCATGA 2448  
+3: L R G H W P A L L \* V L P P P A \*  
+2: T R A L A G S S L S S T A S C M M  
+1: Y E G T G R L F F E F Y R L L H D

DNA: TGC GCGGCCCAAGGAGGGAGATGATCGCCCCCTTCTTCTGGCTCTTTGAGAA 2499  
+3: C A A Q G G R \* S P L L L A L \* E  
+2: R G P R R E M I A P S S G S L R M  
+1: A R P K E G D D R P F F W L F E N

DNA: TGTGGTGGCCATGGGCGTTAGTGACAAGAGGGACATCTCGCGATTTCTCGA 2550  
+3: C G G H G R \* \* Q E G H L A I S R  
+2: W W P W A L V T R G T S R D F S S  
+1: V V A M G V S D K R D I S R F L E

DNA: GTCCAACCCTGTGATGATTGATGCCAAAGAAGTGTCAGCTGCACACAGGGC 2601  
+3: V Q P C D D \* C Q R S V S C T Q G  
+2: P T L \* \* L M P K K C Q L H T G P  
+1: S N P V M I D A K E V S A A H R A

DNA: CCGCTACTTCTGGGGTAACCTTCCCGGTATGAACAGGCCGTTGGCATCCAC 2652  
+3: P L L L G \* P S R Y E Q A V G I H  
+2: A T S G V T F P V \* T G R W H P L  
+1: R Y F W G N L P G M N R P L A S T

DNA: TGTGAATGATAAGCTGGAGCTGCAGGAGTGTCTGGAGCATGGCAGGATAGC 2703  
+3: C E \* \* A G A A G V S G A W Q D S  
+2: \* M I S W S C R S V W S M A G \* P  
+1: V N D K L E L Q E C L E H G R I A

DNA: CAAGTTCAGCAAAGTGAGGACCATTACTACGAGGTCAAACCTCCATAAAGCA 2754  
+3: Q V Q Q S E D H Y Y E V K L H K A  
+2: S S A K \* G P L L R G Q T P \* S R  
+1: K F S K V R T I T T R S N S I K Q

DNA: GGGCAAAGACCAGCATTTTCTGTCTTCATGAATGAGAAAGAGGACATCTT 2805  
+3: G Q R P A F S C L H E \* E R G H L  
+2: A K T S I F L S S \* M R K R T S Y  
+1: G K D Q H F P V F M N E K E D I L

DNA: ATGGTGCACCTGAAATGGAAAGGGTATTTGGTTTCCAGTCCACTATACTGA 2856  
+3: M V H \* N G K G I W F P S P L Y \*  
+2: G A L K W K G Y L V S Q S T I L T  
+1: W C T E M E R V F G F P V H Y T D

DNA: CGTCTCCAACATGAGCCGCTTGGCGAGGCAGAGACTGCTGGGCCGGTCATG 2907  
+3: R L Q H E P L G E A E T A G P V M  
+2: S P T \* A A W R G R D C W A G H G  
+1: V S N M S R L A R Q R L L G R S W

DNA: GAGCGTGCCAGTCATCCGCCACCTCTTCGCTCCGCTGAAGGAGTATTTTGC 2958  
+3: E R A S H P P P L R S A E G V F C  
+2: A C Q S S A T S S L R \* R S I L R  
+1: S V P V I R H L F A P L K E Y F A

DNA: GTGTGTGTAAGGGACATGGGGGCAAACCTGAGGTAGCG 2995  
+3: V C V R D M G A N \* G S  
+2: V C K G H G G K L R \*  
+1: C V \* G T W G Q T E V A

# EXHIBIT C

LETTER

# Estimation of Errors in "Raw" DNA Sequences: A Validation Study

Peter Richterich<sup>1</sup>

Genome Therapeutics Corp., Waltham, Massachusetts 02154 USA

As DNA sequencing is performed more and more in a mass-production-like manner, efficient quality control measures become increasingly important for process control, but so also does the ability to compare different methods and projects. One of the fundamental quality measures in sequencing projects is the position-specific error probability at all bases in each individual sequence. Accurate prediction of base-specific error rates from "raw" sequence data would allow immediate quality control as well as benchmarking different methods and projects while avoiding the inefficiencies and time delays associated with resequencing and assessments after "finishing" a sequence. The program PHRED provides base-specific quality scores that are logarithmically related to error probabilities. This study assessed the accuracy of PHRED's error-rate prediction by analyzing sequencing projects from six different large-scale sequencing laboratories. All projects used four-color fluorescent sequencing, but the sequencing methods used varied widely between the different projects. The results indicate that the error-rate predictions such as those given by PHRED can be highly accurate for a large variety of different sequencing methods as well as over a wide range of sequence quality.

In DNA sequencing, knowledge about the accuracy of sequences can be very valuable. For example, different large-scale sequencing projects may produce sequences at similar rates and costs but with significantly different error rates in the final sequence. One major determinant in the final error rate is the accuracy of the "raw" sequence. Knowledge about the frequency and location of errors in the raw sequence data can help to direct "polishing" efforts to the places where additional effort is needed; it also enables the comparison between different sequencing projects without requiring that the same region be sequenced in each project.

Another area where estimates about sequence error rates would be beneficial is technology development. Accurate error estimates at each base would enable "quality benchmarking" between different methods, thus enabling researchers to choose the method that fills their needs for accuracy and throughput best.

Several groups have developed mathematical models to predict the error probability at any given position in raw sequences. Lawrence and Solovyev used linear discriminant analysis to calculate separate probability estimates for insertions, deletions, and mismatches (Lawrence and Solovyev 1994). Ewing and Green (1998) developed the program

PHRED, which calculates a quality score at each base. This quality score  $q$  is logarithmically linked to the error probability  $p$ :  $q = -10 \times \log_{10}(p)$  (for a discussion of how quality scores are calculated and what the limitations are, see Ewing et al. (1998). When used in combination with sequence assembly and finishing programs that utilize these error estimates, reliable error probabilities promise to increase the accuracy of consensus sequences and to reduce the efforts required in the finishing phase of sequencing projects (Churchill and Waterman 1992; Bonfield and Staden 1995).

To examine the accuracy of probability estimates made by the program PHRED, we compared the actual and predicted error rates for six different cosmid- or BAC-sized projects that were produced by six different large-scale sequencing centers in the United States. All of these six projects used four-color fluorescent sequencing machines; however, the DNA preparation methods, sequencing enzymes, fluorescent dyes and chemistries, and gel lengths varied significantly between the six groups. Table 1 gives an overview of the sequencing projects analyzed. Table 2 lists the different methods used.

## RESULTS

### Error Rate Prediction Accuracy for Six Projects

A comparison of actual and predicted error rates for the six projects in this study is shown in Table 3.

<sup>1</sup>E-MAIL [peter.richterich@genomecorp.com](mailto:peter.richterich@genomecorp.com); FAX (781) 893-9535.

**Table 1. Summary of Data Sets**

Project	Reads	Aligned bases	Average aligned read length
A	455	416,214	915
B	1277	871,230	682
C	1065	603,655	567
D	834	414,595	497
E	1638	1,149,209	702
F	1885	907,796	482
Total	7154	4,362,699	610

The results indicate that PHRED is very successful in identifying bases with low error probabilities. For example, the 1.28 million bases with quality scores of 4–12 (corresponding to error probabilities between 39.8% and 6.3%) contain a total of 187,926 errors. In contrast, the 1.44 million bases with quality scores between 33 and 42 (corresponding to error probabilities between 0.05% and 0.006%) contain only 237 errors, which translates into a 790-fold lower error rate. The trend toward lower error rates can also be observed for each individual project. In most cases, the actual number of errors is close to the predicted error rate. It is also apparent that the actual error rate is typically lower than the predicted error rate.

Both the high overall accuracy and the tendency to slightly overpredict errors are confirmed by statistical analysis, as shown in Table 4. The correlation between predicted and actual error frequencies is excellent for all projects (Spearman correlation coefficient  $>0.89$ ,  $P < 0.0001$ ). Averaged over all projects, the actual error rate is 84.5% of the predicted error rate; the slope of the relation between predicted and actual error rates differs slightly between projects and ranges from 76.6% to 88.4%. To put these differences between projects in relation, it is worthwhile remembering that PHRED quality scores cover a wide dynamic range: The maximum quality score of 51 corresponds to a 50,000-fold lower predicted error rate than the minimum quality score of 4. Even the relative difference between successive quality is larger than the relative difference in the slopes; for example, a quality score of 10 corresponds to an error probability of 10%, whereas a score of 9 corresponds to an error probability of 12.6%.

A different way of looking at the relation between the actual and predicted error rates is shown

in Figure 1. Here, the error rates as a function of the position within all reads in each of the projects, averaged over 50-base windows, is depicted. For all six projects, the predicted error rates are very close to the actual error rates over the entire length of the sequences. Each project has a characteristic distribution of error rates, which differs from each of the other projects. The minimum error rate differs dramatically between projects. The best projects achieve raw error rates of 0.23%–0.36% in the best region of the sequence read, typically from base 150 to 200. The worst project in the data set had an ~10-fold higher error rate of 2.58%.

Toward the end of sequence reads, the error rates increase and start to exceed 10% between bases 300 and 700. In projects that used mainly short gels (e.g., projects D and F), this increase begins sooner, whereas projects that use longer gels show a markedly longer stretch of low error rates (e.g., projects A and B).

Table 5 summarizes key results for the six projects. The first four projects have similar minimum and average error rates. However, the length of the region where the error rate is below 5% differs significantly, from 403 to 682 bases. The project with the shorter low error rate regions contained larger portions of reads generated on short gels, whereas projects A and B were run exclusively on long gels (ABI373 stretch or ABI377 sequencers). Other factors contributing to differences between the first four projects were differences in sequencing chemistries, production scale, and electrophoresis conditions and machines.

Project E and, in particular, project F, had significantly higher error rates than the first four projects. In projects E and F, every sequence generated for the project had been included in the data set, whereas the other four projects had eliminated some "bad" sequences through manual or auto-

**Table 2. Overview of Sequencing Methods Used in the Different Projects**

Template DNA	single-stranded M13, double-stranded plasmids
Sequencing enzymes	Sequenase, Taq, KlenTaqTR, AmpliTaq FS
Sequencing chemistries	Dyes primer (two different dyes chemistries), dye terminator
Sequencing machines	ABI 373, ABI 373 stretch, ABI 377
Gel length	Only short gels, only long gels, mixes of short and long gels

**Table 3. Comparison of Predicted and Actual Error Rates for Six Different Sequencing Projects**

Project	Quality score	4-12	13-22	23-32	33-42	43-51
A	aligned bases	119,246	75,293	70,391	144,876	73,234
	expected errors	20,256	2,064	172	37	1
	actual errors	16,784	1,758	127	17	1
B	aligned bases	182,034	137,940	181,998	399,690	140,176
	expected errors	29,953	3,704	410	102	3
	actual errors	26,038	2,536	287	35	0
C	aligned bases	139,345	131,419	151,197	292,070	68,529
	expected errors	22,277	3,411	357	74	2
	actual errors	16,670	1,513	194	26	3
D	aligned bases	103,898	68,995	68,613	153,730	111,752
	expected errors	16,880	1,919	168	38	3
	actual errors	14,495	1,924	146	59	2
E	aligned bases	378,755	217,438	167,968	392,717	144,313
	expected errors	63,947	6,336	418	95	4
	actual errors	55,968	6,516	355	67	5
F	aligned bases	359,809	136,688	98,840	64,035	5,130
	expected errors	66,938	4,079	256	23	0
	actual errors	57,971	3,856	332	33	1
All	aligned bases	1,283,087	767,773	739,007	1,447,118	543,134
	expected errors	220,252	21,513	1,781	370	13
	actual errors	187,926	18,103	1,441	237	12

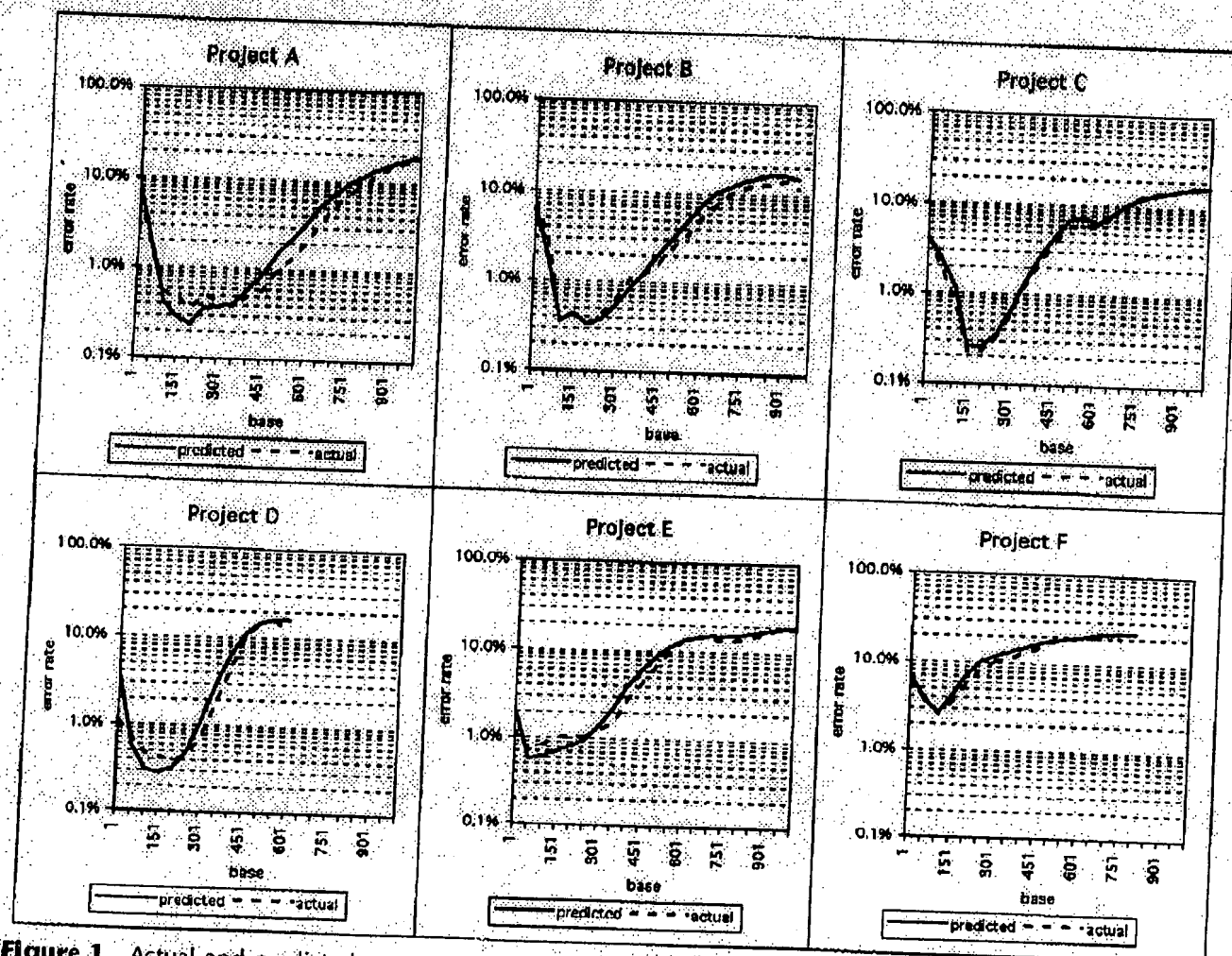
matic inspection. After eliminating <10% of the worst sequences in project E, the error rates for the remaining sequences were comparable to those of the first four projects. In contrast, project F showed a much more uniform distribution of sequence quality.

The last column in Table 5 shows the average number of bases with an estimated error probability of at most 0.1%, which is equivalent to a quality score of at least 30. The count of such "very high-quality" bases is a good indicator of sequence quality, both for individual sequences and, when aver-

**Table 4. Summary of Statistical Analysis Results**

Project	Spearman p	$P >  p $	Slope	t ratio	$P >  t $
A	0.9646	<0.0001	0.818	75.1	<0.0001
B	0.9890	<0.0001	0.874	98.2	<0.0001
C	0.9846	<0.0001	0.766	71.6	<0.0001
D <sup>a</sup>	0.8692	<0.0001	0.855	68.3	<0.0001
E	0.9956	<0.0001	0.884	144.3	<0.0001
F	0.9968	<0.0001	0.865	151.6	<0.0001
All	0.9964	<0.0001	0.845	174.5	<0.0001

<sup>a</sup>In project D, the Spearman correlation coefficient p was artificially low as only very few bases (10) bases had a quality score of 5, and none of these bases contained an actual error (expected: 3.16 errors). Exclusion of this quality score gave a Spearman correlation coefficient of 0.9786 ( $P < 0.0001$ ). The frequencies in the slope calculations were weighed by the number of bases at any given quality score and, thus, were not sensitive to such small sample distortions (see Methods).



**Figure 1** Actual and predicted error rates in six different sequencing projects. Actual error rates and predicted error rates in 50-base windows over the length of the sequence reads, averaged over all reads that could be aligned to the consensus sequence by CROSS\_MATCH, are shown. The numbers on the x-axis show the first base in a given 50-base window.

aged over all sequences in a project, as an indicator for the entire project. Compared to the estimated error rates, the count of very high-quality bases is less prone to distortions from a small number of low-quality reads, as the data for project E demonstrate.

#### Prediction Accuracy for Data Subsets of Different Quality

The quality of sequences within any given project can vary substantially, and the use of predicted error rates has the potential to be a powerful tool for qual-

**Table 5. Comparison of Key Results for Six Different Sequencing Projects**

Project	Actual minimum error rate (%)	Actual average error rate (%)	Length of <1% error region	Length of <5% error region	Average bases with $P(\text{error}) < 0.1\%$
A	0.36	3.6	422	682	468
B	0.34	2.8	274	567	395
C	0.23	2.4	291	479	348
D	0.39	3.1	300	403	294
E	0.71	4.7	129	464	317
F	2.58	9.2	0	162	79



ity analysis and control in large-scale DNA sequencing projects. To analyze how accurate PHRED error estimates are for different quality sequences within the same sequencing project, we subdivided a data set into four quartiles, based on the number of very high-quality bases in each sequence (see Methods). The comparison of actual and predicted error rates is shown in Figure 2.

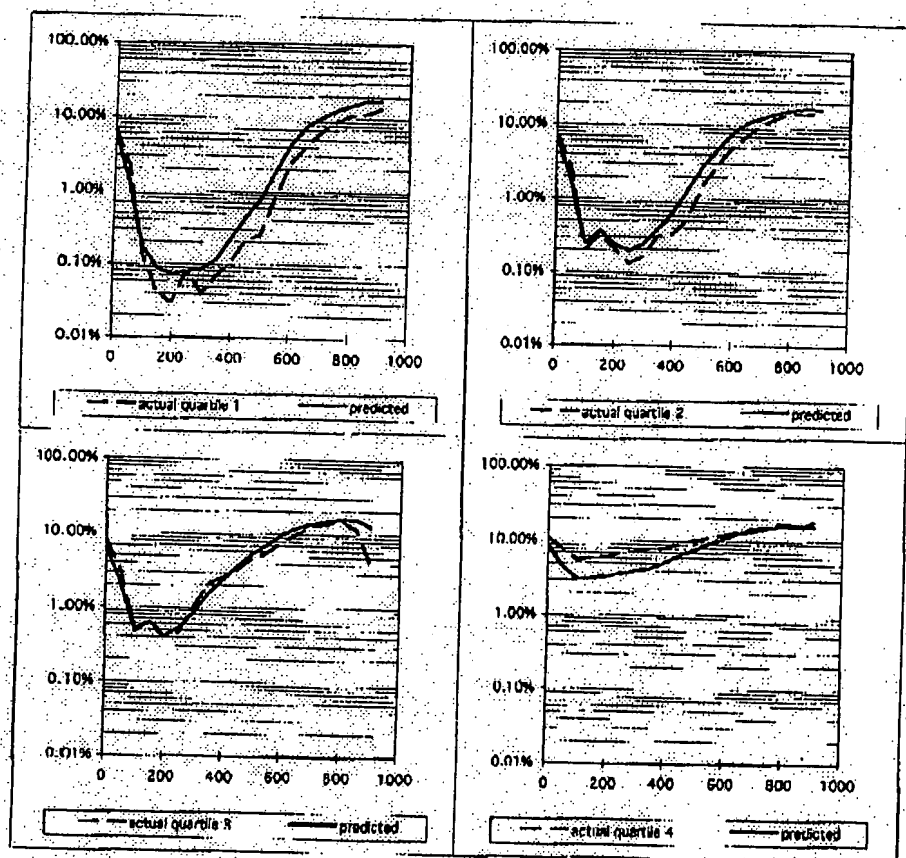
When measured by the error rate in the best region of a sequence, the data quality in the different quartiles varies >100-fold between the best and the worst 25% of the sequences. The best quartile showed ~0.03% error for >100 bases, whereas the error rate in the worst quartile always exceeded 5%. In quartiles 2 and 3, the predicted error rates match the actual error rates very closely. In the best and

worst quartiles, PHRED's accuracy was somewhat lower from base 100 to 500. In the best sequences, PHRED's error estimates were about twofold too high; in the worst sequences, the error estimates were too low, again by a factor of 2. This underprediction of errors can be partially explained by the fact that PHRED gives ambiguous base calls (N's) a quality score of 4, corresponding to an error probability of 39.8%; however, N's will always show up as an actual error. Even in the worst and best quartiles, however, the predicted error rate curves are very similar to the actual error rate curves.

The results shown in Figure 2 also demonstrate that the count of very high-quality bases, or bases with an estimated error probability of at most 0.1%, can be used effectively to characterize the overall

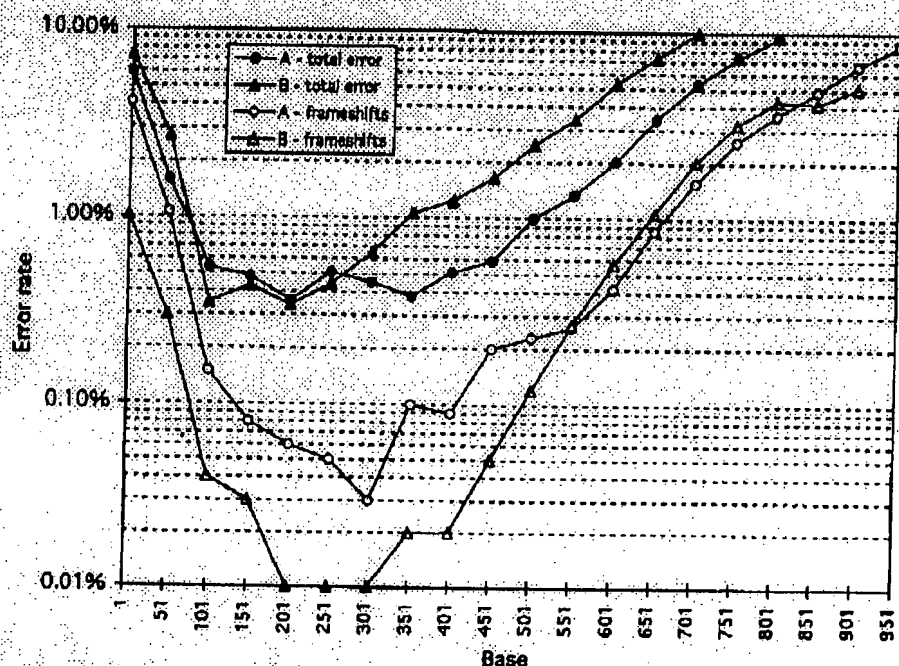
quality of a sequence read. Sorting the sequence reads into quartiles based on the number of very high-quality bases worked well, as shown by the >100-fold difference in the minimum error rate between the first and the fourth quartile.

Other methods to characterize the overall quality of individual reads based on PHRED quality scores can give similar results. For example, counting bases above a minimum quality threshold anywhere in the range of 20–40 gave similar results for most data sets (not shown), and such counts are used by a number of different laboratories as quality measures. Alternatively, the quality values can be converted to error probabilities and averaged to give the predicted error rate for the trace, or summed to give the total predicted number of errors in a trace. However, such averages and totals can sometimes give a misleading picture, as the following example illustrates. Assume that two sequence reads have very similar quality in the alignable part of the read but that one of the two sequences was run much longer and



**Figure 2** Actual and predicted error rates in different quality subsets of project B. Sequence reads were sorted by the number of bases with a predicted error rate of at most 0.1% (very high-quality bases), and assigned to quartiles, with quartile 1 corresponding to the highest numbers. Actual and predicted error rates for all sequences in each subset were calculated as in Fig. 1. Note that a number of sequence reads that had been rejected because of too low quality were added back to the data set for illustrative purposes, all of which are in quartile 4. These sequences were not included in the data sets used to generate Figs. 1 and 3 and Tables 1 and 3.

## RICHTERICH



**Figure 3** Actual frameshift and total error rates for projects A and B. To calculate frameshift error rates, only insertions and deletions were counted. Mismatch errors, which account for the vast majority of errors after base 150, were included only in the total error count. Note that project B ( $\Delta$ ,  $\triangle$ ) has a slightly similar or slightly higher total error rate compared to project A ( $\bullet$ ,  $\circ$ ) but only about one-third as many insertions and deletions up to base 500. For both projects, the frameshift error rate in the raw data is  $<1$  in 1000 for  $>300$  bases, and  $<1$  in 10,000 for  $>100$  bases in project B.

therefore contains a longer unalignable "tail" of very low-quality bases. When calculating the average error rate for these two sequences, the second sequence will have a much higher average error and, therefore, appear to be of lower quality. In contrast, the counts of very high-quality bases for both sequences will be very similar, as the unalignable tails contain few, if any, high-quality bases. Therefore, counts of bases above a high enough quality threshold will give a more robust and clearer picture of trace quality.

### Frameshift Error Rates for Different Sequencing Chemistries

Depending on how biologists use DNA sequences, knowledge about total error rates in raw sequences may or may not be sufficient. For example, frameshift errors in coding sequences will generally lead to incorrectly predicted open reading frame, whereas mismatch errors will do so only if the mismatch introduces a stop codon or a new splice site. At the time of this writing, PHRED did not differentiate between mismatch and frameshift errors, but only estimated total error rates. This might occa-

sionally lead to questionable conclusions, as the results shown in Figure 3 illustrate.

Figure 3 shows the total actual error rates and the frameshift error rates for two projects, A and B. The total error rates for both projects are similar for up to 350 bases; after 350 bases, project B has a somewhat higher total error rate. However, examining the frameshift error rate gives rise to a different picture: from base 1 to 500, project A has approximately four times as many insertions and deletions as project B. This difference in frameshift error rates can be explained by the sequencing chemistries that were used in the two projects. Project B, with the lower frameshift error rate, used only dye terminator chemistry, which is known to eliminate band spacing artifacts from hairpin structures ("compressions"). Project A, on the

other hand, used dye primer chemistry, which is more prone to insertion and deletion errors from mobility artifacts, for most sequencing reactions.

## DISCUSSION

As large-scale DNA sequencing has become a more routine and common process, the traditional methods for assessing sequence quality have become unsatisfactory. In projects like single-pass cDNA sequencing, it is not possible to calculate and compare error rates after finishing a sequence, as finishing never takes place. Even when a comparison between raw and finished sequence can be done, the time delay between raw data generation and quality assessment is often large. This delay makes it difficult to improve ongoing projects, and it sometimes makes it impossible to capture problems early on. Some immediate quality feedback can be reached by including known standard sequences for quality control. However, this approach can be costly, and it fails when error profiles differ between standard and unknown sequences.

In contrast to these traditional methods to assess sequence accuracy, direct estimation of error

rates in raw sequence data would enable immediate quality control and feedback. Accurate, base-by-base estimates of error probabilities could also increase the utility of single-pass sequences significantly, allow efficient comparison and optimization of different sequence chemistries, and enable the development of better software tools for sequence assembly and analysis.

The critical question for any error rate prediction tool is how accurate are the error rate estimates, in particular if different sequencing methods and chemistries are used? The results presented herein provide an answer to this question for the program PHRED, as well as clues where further development would be useful. As shown in Tables 3 and 4 and in Figure 1, the agreement between predicted and actual error rates was very good in each of the six different projects analyzed. The observed high level of prediction accuracy in all of these projects is almost astonishing if one takes into account that actual errors are binary (a base is either correct or wrong), whereas predicted error rates are probabilities on a scale from 0.0 to 1.0. The observed tendency to overpredict error rates can be at least partially explained by the "small sample correction" that was used in the derivation of threshold parameters for quality scores (Iwling and Green 1998). For most practical applications, such a somewhat conservative estimation of quality scores is tolerable or even desirable. Overall, the results clearly show that error probabilities given by PHRED accurately describe raw sequence data quality.

In judging the usefulness of predicted error probabilities, it is important to know how differences in sequencing methods will influence the prediction accuracy. For example, the larger variation in peak heights tends to be larger in dye terminator sequencing than in dye primer sequencing, and different sequencing enzymes are known to produce different specific height variation patterns. Any estimation of error probabilities that takes the peculiarities of a specific sequencing chemistry into account would therefore be expected to be less accurate for different chemistries.

The projects included in this study were specifically chosen to provide an initial answer to the question of how generally useful PHRED quality scores are. These projects represent the vast majority of different multicolor fluorescent sequencing methods used in the last 3 years: different template DNAs and DNA preparation methods, different enzymes, gel lengths, run conditions, and different fluorescent dyes. The data also include a considerable spread in data quality, both between projects

and within individual projects. None of the projects analyzed here were included in PHRED's training set, and just one of the six laboratories that contributed data to this study also contributed data to the training data sets. One of the projects in this study consisted entirely of dye terminator sequences, which presented only a small fraction of the sequences in the test data set. Another project exclusively used a set of fluorescent dyes different from those used in the training sets. Each project differed from the other projects in this study in at least one, and typically many, experimental aspects like template preparation, sequencing enzymes, gel run conditions, and so forth. Despite these differences, the accuracy of error rate predictions was very similar for all projects.

Our results justify some optimism about the accuracy of PHRED quality scores for minor changes in sequencing technology, for example, sequences generated by new enzymes and fluorescent dyes. Initial studies showed that PHRED quality scores were also accurate for sequences produced by multiplex sequencing with radioactive detection (P. Richterich, unpubl.). However, we also observed two effects that can invalidate PHRED quality scores during these studies. First, sequences generated by chemical sequencing gave too low quality scores at mixed (A + G) reactions. Because secondary peak height is one of the parameters used in the error rate predictions, this is not surprising. Another potential source of error is high-frequency noise in the trace data. With such data, PHRED occasionally underestimated the band spacing by a factor of 2 or more, which resulted in incorrect base calls and quality scores. By applying simple smoothing algorithms to data with high-frequency noise, these problems could typically be resolved. Similar steps may be necessary to obtain accurate PHRED quality scores on data that have been generated by different sequencing instruments or preprocessed by different software.

Accurate quality scores can have a major impact on how sequences are used downstream from the sequence production process. In traditional sequencing projects where the goal is complete coverage at a final error rate below (e.g.) 1 in 10,000, the accuracy goals can be reached with single sequence reads as long as the quality scores are at least 40 (however, other potential problems like clone instability may make higher coverage advisable). Interesting questions arise as to how individual read quality contributes to project quality, or the error rate of the "final" sequence. Under the assumption that errors between different sequence reads are

## RICHTERICH

completely independent, one could argue that two reads with a quality score of 20 (error probability of 1 in 100) are just as valuable as one sequence with a quality score of 40 (error probability of 1 in 10,000). However, although a single sequence stretch with quality levels above 40 would give a final sequence with an error rate of <1 in 10,000, assembling a consensus from two sequences with quality scores of 20 (1% error rate) could lead to one of two results: If the errors were completely random, the consensus sequence would be ambiguous at 2% of all locations; if the errors were completely localized, for example, because of reproducible compressions, the consensus sequence would have one "hidden" error every 100 bases. Typically, consensus sequences derived from low-quality sequences will have both kinds of problematic regions. Increased coverage can rapidly eliminate the random errors; however, increased coverage does not resolve errors from systematic sources. Manual examination of such problem areas is generally required; such "contig editing," however, tends to be time consuming, requires highly trained personnel, is an obstacle toward complete automation of DNA sequencing, and sometimes fails to eliminate all errors. This leads to the somewhat counterintuitive conclusion that the practical value of increasing sequence quality can be even higher than indicated by the quality scores: One sequence of average quality above 40 can be "worth" more than two sequences of average quality 20.

Another application of DNA sequencing where high quality can be of disproportionately high value is the search for mutations in genomic DNA. In low quality sequences, secondary peaks and low resolution often complicate the identification of heterozygous mutations. In regions of higher sequence quality, such secondary peaks are smaller or absent and peaks are better resolved. Therefore, both false-positive and false-negative errors can be significantly reduced in high-quality regions. Tools like PHRED, which can accurately measure sequence quality from trace data, can be of twofold value for mutation detection. First, base-specific quality scores can allow optimization of sequencing methods and strategies for mutation detection. Second, the quality scores can be used to evaluate the usefulness of individual sequence reads for mutation detection (e.g., by discarding reads below minimum thresholds), and they can guide software that automatically detects mutations.

The ability to predict error rates in a highly accurate fashion is likely to have a major impact in applications like those described above. PHRED is

the first widely used program that accurately predicts base-specific error probabilities. However, the algorithm for determining quality values has been described (Ewing and Green 1998), and it should be straightforward to implement similar quality values in other base-calling programs. Furthermore, an extension of the approach developed by Ewing and Green should be possible. For example, differentiation between mismatch and frameshift errors would enable better comparisons of sequencing methods with similar total error rates but different frameshift error rates. Several groups have described efforts to calculate separate probabilities (or "confidence assessments") for mismatch errors and frameshift errors (Lawrence and Solov'yev 1994; Berno 1996). Their results demonstrated that different approaches to error type characterization are feasible and promising. Implementation of such error type predictions in other programs similar to the way PHRED uses quality scores would enable better method assessments, benchmarking, and production quality control, and could have a significant impact on downstream uses of DNA sequence information.

## METHODS

### Data Sets

For one project, sequence raw data in the form of ABI trace files were downloaded from a public FTP site. Sequence data for the five other projects were kindly provided by five different large-scale sequencing groups. Table 1 gives a summary of the six projects, and Table 2 gives an overview of the different sequencing methods used in the projects. The projects differed in the amount of prescreening of data that had been done, reflecting different approaches to quality control in different laboratories. In two projects (B and C), different software programs had been used to identify and eliminate low-quality sequences. One project (F) included all data files generated, whereas the other three projects had excluded "failed lanes."

### Comparison of Actual and Predicted Error Rates

The sequences for all traces in each project were recalled using the program PHRED (v. 961028). Next, sequences in each project were assembled with PHRAP (P. Green, unpubl.). Slightly different methods were chosen for the statistical and graphical evaluation of the error rate prediction accuracy. In the statistical evaluation, only the longest contig produced by PHRAP was considered. The tables of aligned bases and observed discrepancy counts for



each quality score were taken from the PHRAP output and analyzed as follows. The expected number of discrepancies ( $E$ ) at each quality score ( $q$ ) was calculated by multiplying the number of aligned bases ( $N$ ) with the error probability corresponding to the quality score:  $E = N \cdot 10^{-q/10}$ . The Spearman ranking coefficients were calculated by comparing the expected and observed error frequencies. To obtain the quantitative relation between the expected and observed error rates over the entire range, a least-squares fit between the observed and expected rates was performed, with the intercept set to zero and the number of aligned bases at each quality score used as weights.

For a graphical comparison of estimated and actual error rates in 50-bp windows, the following steps were taken. For two of the projects, the consensus sequence was retrieved from public databases. For the four other projects, the DNA sequence and quality information were used by the program PHRAP to assemble consensus sequences for each of the projects. The individual reads were aligned to the consensus sequences of the longest contig, using the program CROSS\_MATCH (P. Green, unpubl.), after removing single-coverage regions from the ends of the consensus sequence. CROSS\_MATCH uses an implementation of the Smith-Waterman algorithm to generate alignments that typically do not include the ends of sequences, where disagreements are commonly due to vector sequence or low quality sequence.

The quality files generated by PHRED and the alignment summaries generated by CROSS\_MATCH were then analyzed as follows. First, the region of each query sequence that had been aligned by CROSS\_MATCH was determined. Next, the actual and predicted error rates for the entire aligned part of each individual sequence was calculated. In addition, the average actual and predicted error rates for all alignable sequences together were calculated for windows of 50 bases in length. To calculate the predicted error rate, the quality scores  $q$  determined by PHRED at each base were converted to error probabilities as described above (Ewing and Green 1998).

#### Subdividing Data into Subsets Based on Data Quality

To examine the accuracy of PHRED quality scores for data subsets of different quality within a project, the following approach was taken. For all sequence reads in project B, the number of bases with a quality score of at least 30 in each sequence was determined (bases with quality scores of at least 30 were called very high-quality bases, or VHQ bases). Se-

quences were sorted in descending order based on the number of very high-quality bases, and divided into four quartiles. Accordingly, quartile 1 contained 25% of sequences with the highest number of very high-quality bases, and quartile 4 contained the "worst" sequences. To illustrate the prediction accuracy in data with relatively high error rates, sequences from project B that had been "discarded" because they had not met the minimum quality criteria were added back to the data set. The sequences in each quartile were compared to the consensus sequences that had been generated using the entire data set, as described above for the graphical comparison.

#### Determining Actual Frameshift Error Rates

The calculation of actual frameshift error rates in the raw sequence data was performed using CROSS\_MATCH, similar to the procedure described above for total error rates, except that only insertion and deletion errors were counted. Because PHRED does not give separate frameshift error estimates, a comparison of predicted and actual frameshift errors is not possible.

#### ACKNOWLEDGMENTS

I thank the participating laboratories for contributing their data, Dr. Josée Dupuis for help with the statistical analysis, and Dr. Phil Green for helpful discussions.

The publication costs of this article were defrayed in part by payment of page charges. This article must therefore be hereby marked "advertisement" in accordance with 18 USC section 1734 solely to indicate this fact.

#### REFERENCES

- Bemo, A.J. 1996. A graph theoretic approach to the analysis of DNA sequencing data. *Genome Res.* 6: 80-91.
- Bonfield, J.K. and R. Staden. 1995. The application of numerical estimates of base calling accuracy to DNA sequencing projects. *Nucleic Acids Res.* 23: 1406-1410.
- Churchill, G. and M.S. Waterman. 1992. The accuracy of DNA sequences: estimating sequence quality. *Genomics* 14: 89-98.
- Ewing, B. and P. Green. 1998. Base-calling of automated sequencer traces using phred. II. Error probabilities. *Genome Res.* (this issue).
- Ewing, B., L. Hillier, M.C. Wendt, and P. Green. 1998. Base-calling of automated sequencer traces using phred. I. Accuracy assessment. *Genome Res.* (this issue).
- Lawrence, C.B. and V.V. Soloviyev. 1994. Assignment of position-specific error probability to primary sequence data. *Nucleic Acids Res.* 22: 1272-1280.

Received October 27, 1997; accepted in revised form February 3, 1998.

# GENOME RESEARCH

Volume 8 Number 3  
March 1998

## Editors

**Laurie Goodman**  
Cold Spring Harbor Laboratory  
**Mark Boguski**  
National Center for Biotechnology  
Information, NIH  
**Aravinda Chakravarti**  
Case Western Reserve University

**Richard Gibbs**  
Baylor College of Medicine  
**Eric Green**  
National Human Genome  
Research Institute, NIH  
**Richard Myers**  
Stanford University School of Medicine

## Editorial Board

**Rakesh Anand**  
Zeneca Pharmaceuticals  
**Spyros Antimarakis**  
University of Geneva  
**Charles Auffray**  
CNRS  
**Philip Avner**  
Institut Pasteur  
**Andrew Ballabio**  
Telethon Institute of Genetics and  
Medicine  
**David Bentley**  
The Sanger Centre  
**Bruce Birren**  
Whitehead Institute/MIT Center for  
Genome Research  
**Michael Boehnke**  
University of Michigan School of  
Public Health  
**Annie Bowcock**  
University of Texas Southwestern  
Medical Center  
**David Burke**  
University of Michigan Medical School  
**Jeffrey Chamberlain**  
University of Michigan Medical School  
**Ellson Chen**  
Perkin-Elmer Corporation  
**David R. Cox**  
Stanford University School of Medicine  
**Ronald W. Davis**  
Stanford University School of Medicine  
**Richard Durbin**  
Sanger Centre, UK  
**Joseph Ecker**  
University of Pennsylvania  
**Devery S. Emanuel**  
Children's Hospital of Philadelphia  
**Raymond Henwick**  
Biotech Laboratories  
**Chris Fields**  
National Center for Genome Resources  
**Simon Foote**  
Walter and Eliza Hall Institute of  
Medical Research

**Phil Green**  
University of Washington  
**Kenshi Hayashi**  
Kyushu University  
**Philip Hieter**  
The Johns Hopkins University School  
of Medicine  
**Clare Huxley**  
St. Mary's Hospital Medical School  
**Howard J. Jacob**  
Medical College of Wisconsin  
**Alec Jeffreys**  
University of Leicester  
**Mark Johnston**  
Washington University School of  
Medicine  
**Mary-Claire King**  
University of Washington  
**Ben Koop**  
University of Victoria  
**Pui-Yan Kwok**  
Washington University School of  
Medicine  
**Ulf Landegren**  
Uppsala Biomedical Center  
**Mark Lathrop**  
The Wellcome Trust Centre  
**Michael Lovett**  
University of Texas Southwestern  
Medical Center  
**Jen-I Mun**  
Genome Therapeutics Corporation  
**Douglas Marchuk**  
Duke University Medical Center  
**Thomas Marr**  
Cold Spring Harbor Laboratory  
**W. Richard McCombie**  
Cold Spring Harbor Laboratory  
**Susan Naylor**  
University of Texas Health Science  
Center  
**David Nelson**  
Baylor College of Medicine

## Reviews Editor

**Allison Stewart**  
Cambridge, UK

STEENBOCK  
MEMORIAL LIBRARY

APR 13 1998

U.W.-MADISON

**Maynard Olson**  
University of Washington  
**Svante Pääbo**  
University of Munich  
**Leena Peltonen**  
National Public Health Institute,  
Helsinki  
**David Porteous**  
MRC Human Genetics Unit  
Western General Hospital, Edinburgh  
**Roger Reeves**  
Johns Hopkins University School of  
Medicine  
**Bruce Roe**  
University of Oklahoma  
**Rodney Rothstein**  
Columbia University College of P&S  
**Gerald Rubin**  
University of California, Berkeley  
**Lloyd Smith**  
University of Wisconsin-Madison  
**Randall Smith**  
Baylor College of Medicine  
**Marcelo Bento Soares**  
University of Iowa  
**William Studier**  
Brookhaven National Laboratory  
**Grant Sutherland**  
Women's and Children's Hospital,  
Adelaide  
**Barbara Trask**  
University of Washington  
**Gert-Jan B. van Ommen**  
Leiden University  
**Robert D. Wells**  
University of Utah  
**Jean Weissenbach**  
Genethon, CNRS  
**Richard Wilson**  
Washington University School of  
Medicine  
**James Womack**  
Texas A&M University

## Editorial Office

Cold Spring Harbor Laboratory Press  
1 Bungtown Road  
Cold Spring Harbor, New York 11724  
Phone (516) 367-6834  
Fax (516) 367-6334  
<http://www.cshl.org>

## Editorial/Production

**Nadine Dumser**, Technical Editor  
**Kristin Kraus**, Production Editor  
**Cynthia Grimm**, Production Editor  
**Peggy Callicchia**, Editorial Secretary